

Knowledge Transfer across Imaging Modalities Via Simultaneous Learning of Adaptive Autoencoders for High-Fidelity Mobile Robot Vision

Md Mahmudur Rahman¹, Tauhidur Rahman², Donghyun Kim², and Mohammad Arif Ul Alam¹

Abstract—Enabling mobile robots for solving challenging and diverse shape, texture, and motion related tasks with high fidelity vision requires the integration of novel multimodal imaging sensors and advanced fusion techniques. However, it is associated with high cost, power, hardware modification, and computing requirements which limit its scalability. In this paper, we propose a novel Simultaneously Learned Auto Encoder Domain Adaptation (SAEDA)-based transfer learning technique to empower noisy sensing with advanced sensor suite capabilities. In this regard, SAEDA trains both source and target auto-encoders together on a single graph to obtain the domain invariant feature space between the source and target domains on simultaneously collected data. Then, it uses the domain invariant feature space to transfer knowledge between different signal modalities. The evaluation has been done on two collected datasets (LiDAR and Radar) and one existing dataset (LiDAR, Radar and Video) which provides a significant improvement in quadruped robot-based classification (home floor and human activity recognition) and regression (surface roughness estimation) problems. We also integrate our sensor suite and SAEDA framework on two real-time systems (vacuum cleaning and Mini-Cheetah quadruped robots) for studying the feasibility and usability.

Index Terms—Radar, Lidar, Domain Adaptation, Activity Recognition, Surface Characterization, Quadruped robot

I. INTRODUCTION

Recent advancements of mobile robotic vision technologies (such as LiDAR, Camera, Depth, Infrared, Radar, Terahertz sensors) enable quadruped robots to navigate walking paths for blinds [1], play with owner [2], guardhouse from intruders [3], remote monitoring of oil, gas and power installations and the construction sites [4]. However, all of the above applications are involved a task-specific sensor suite which limits the scalability of the system [1]–[4]. Some of the sensors are more powerful than others in the task-specific applications, such as Lidars are generally found much superior to Radar in empowering mobile robots to solve versatile shape and movement sensing problems [5]. Recent advancements in radar imaging technology hold significant promise in versatile sensing applications with mobile robots due to affordability, reliability, and higher material penetrability (surface properties). But, radar still falls short in capturing precise and high-fidelity images of objects,

surfaces, and human activities with a signal-to-noise ratio as high as Lidar [5]. To accommodate the capabilities of both of the sensors, recent quadruped robots use both imaging modalities and sometimes more advanced ones along with novel fusion techniques and algorithms facilitating significant improvement in robotic vision-based solutions [6]–[8]. Thus, scaling capabilities of existing mobile robots have become impossible without modifying existing hardware and algorithms, which forces customers to replace the old system with the latest one. In this paper, we present a way of improving the capabilities of existing quadruped robot vision, by utilizing a temporarily designed advanced sensor suite to align with existing sensors, collect few simultaneous data and use that to train a new algorithm to sustain the advanced sensor suite’s capabilities without its presence i.e., empowering an existing mobile robot’s capabilities without any permanent hardware modifications.

Several recent works have explored simultaneous learning and knowledge transfer among multiple sensing modalities in solving complex robotic autonomy, navigation, and vision problems. Weston et. al. [9] proposed a self-supervised deep learning-based method to estimate Radar sensor grid cell occupancy probabilities considering simultaneously collected Lidar-based probabilities as labels to train a deep regression model. [10] developed a deep learning framework for Radar-based multi-object localization where the labels come from the sophisticated Lidar and RGB camera fusion-based models. In [11], researchers show that a Lidar-based localization framework ScanContext [12] can be used with radar data to improve the radar accuracy. An indoor localization using geometric structure was shown in [13] marking the radar points on lidar generated maps and CAD models. Yin et. al. [14] proposed a conditional generative adversarial network (GAN) based domain adaptation technique, the closest work to SAEDA to generate Lidar representations of Frequency-Modulated Continuous-Wave (FMCW) Radar data and learn from it. A Monte Carlo localization (MCL) system is also formulated by motion and measurement models [14]. None of the above methods propose to utilize domain adaptation to enhance robotic vision capabilities aided by a temporarily advanced sensor suite and simultaneous learning schemes on temporarily collected simultaneous data.

In this paper, we aim to answer two **key questions**: (q1) *Can a temporarily integrated sensor suite collected high-dimensional heterogeneous simultaneous data be utilized to*

¹University of Massachusetts Lowell
mdmahmudur_rahman@student.uml.edu,
mohammadariful_alam@uml.edu

²University of Massachusetts Amherst
trahman@cs.umass.edu, donghyunkim@cs.umass.edu



Fig. 1: Our Lidar-Radar sensor suite integrated with Vacuum Cleaning Robot and Mini-Cheetah quadruped robot [25]

empower existing mobile robot vision with transfer learning? and *(q2) Can we optimize the number of simultaneously collected labeled data i.e., semi-supervised learning?* Recent advancement of deep domain adaptation facilitates significant improvement in target domain classification performance in presence of label scarcity [15], noise [16], and heterogeneity [15] via generative networks ([16]–[18]) or Discrepancy-based methods ([19]–[22]). However, the generative networks-based domain adaptations training process is complex and does not guarantee convergence in the case of high-dimensional signals (Lidar/Radar) *(q1)* [23]. On the other hand, discrepancy-based methods related errors on the target domain are bounded by distribution divergence which makes it difficult to learn a domain invariant feature space between the heterogeneous source and target domain *(q1)* [24]. Finally, all of the existing domain adaptation methods need an extensive amount of labeled source data that challenges *q2* requiring to design of a new domain adaptation algorithm for our target problem.

In this paper, we propose *SAEDA*, an auto-encoder based domain adaptation method which is both easy to train like discrepancy based methods [20], [26], [27] and achieve a state of the art performance like adversarial networks [15], [28] both in semi-supervised use cases [20], [26] in presence of a small amount of labeled target dataset. More specifically, our **key contributions** are as follows:

- We propose a novel simultaneous learning scheme for high dimensional imaging modalities (LiDAR/Radar) via a novel loss function which helps the target feature space to be adapted along the way of learning of the source domain feature space with only simultaneously collected data from existing imaging modality and temporarily designed high fidelity sensor suite.
- Additionally, we design a classification method by utilizing learned domain invariant feature spaces from the source and transferring them to the target to sustain characteristics of both discrepancy and adversarial based domain adaptation characteristics altogether for the semi-supervised scenario.

- We evaluate the performance of *SAEDA* with real-time data on three use cases: (i) human activity recognition for robot-human interaction research, (ii) surface characterization to assist mobile robot for automated vacuum cleaning planning, and (iii) sandpaper surface estimation.
- Finally, we integrate our high fidelity sensor suite (LiDAR+Radar) in the Mini-Cheetah quadruped robot and study the efficacy/feasibility of our framework and system for outdoor surveillance.

II. DEEP DOMAIN ADAPTATION MODELING

A. Setting

We denote the source domain data having n_s numbers of samples as $\mathcal{D}_s = \{(X_s^i, Y_s^i)\}^{n_s}$ where X_s^i and Y_s^i are the features and the class label of the i^{th} sample respectively. Here, $X_s^i \in \mathbb{R}^{d_s}$; d_s is the dimension of the source features. We define the source domain classification task, T_s is to classify the source domain data correctly.

In case of the target domain data-set, we divide it into two dis-join sets, one for labeled (\mathcal{D}_l) and another for unlabeled (\mathcal{D}_u) samples. So, the target domain, $\mathcal{D}_t = \mathcal{D}_l \cup \mathcal{D}_u = \{(X_{tl}^i, Y_{tl}^i)\}^{n_l} \cup \{(X_{tu}^i)\}^{n_u}$. Here, $X_{tl}^i, X_{tu}^i \in \mathbb{R}^{d_t}$ are the labeled and the unlabeled i^{th} target samples respectively where d_t is the dimension of the target domain features. We assume $n_s \gg n_l$ and $n_u \gg n_l$, the number of target labeled sample is very smaller than number of target labeled and source labeled sample in our domain adaptation setup. Similar to T_s , the target classification task T_t is defined to classify the unlabeled target data X_{tu}^i .

B. Adaptive Autoencoder via Simultaneous Learning

Fig. 2 (a) shows the complete training process of *SAEDA* where learning happens in three steps, simultaneous training of the auto-encoders, training the classifier, and fine-tuning.

1) *Training the Auto-encoders via Simultaneous Learning:* In this stage, the source and the target auto-encoders are trained with respective batches of feature vectors simultaneously. Eventually, the encoder parts of the auto-encoder, $E^s(\cdot)$ and $E^t(\cdot)$ map the input features X_s and X_t to the bottleneck feature representation space \hat{X}_s and \hat{X}_t respectively.

$$\hat{X}_s = E^s(X_s), \quad \hat{X}_t = E^t(X_t) \quad (1)$$

During the training, the objective function is set in a way that minimizes the statistical distance between \hat{X}_s and \hat{X}_t . We propose to modify the Maximum Mean Discrepancy (MMD) [29] loss function and introduce the Class-wise MMD which minimizes the domain discrepancy according to the classes along the training. However, we apply a simultaneous learning scheme where both of the source and target auto-encoders get optimized simultaneously on a single graph developed on the top of both auto-encoders. This ensures matching of complex features between the source and the target domain along the way of convergence. We

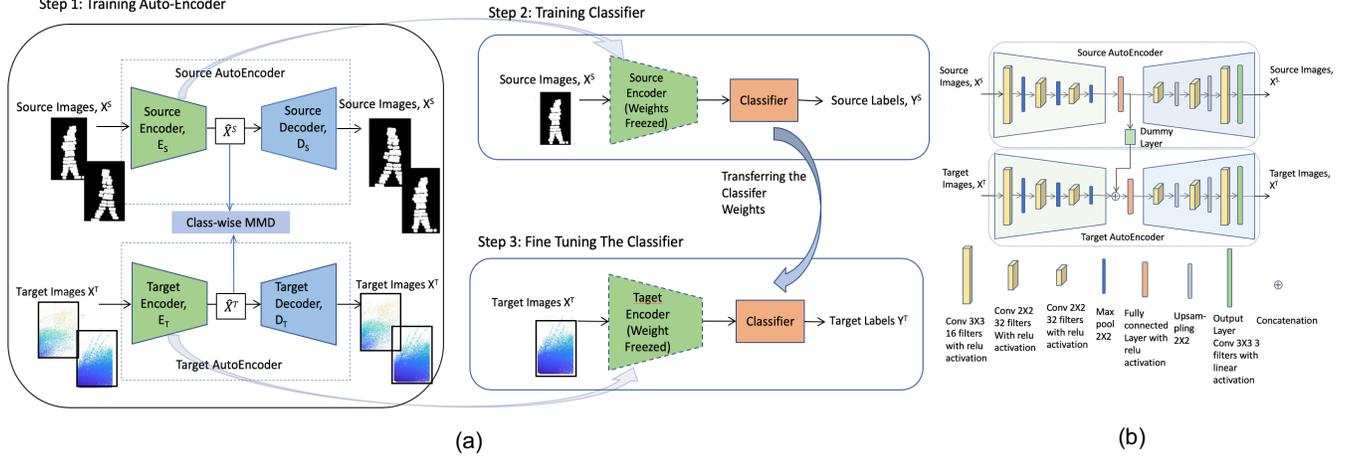


Fig. 2: (a) Training procedures of our SAEDA architecture. The dashed line represents the module is already trained and the weights are frozen in this step. The solid line represents the weights of this module are being updated in this step. (b) Network architecture of the source and the target auto-encoder.

design a self reconstruction loss which sets the source auto-encoder optimization loss as follows.

$$\mathcal{L}_s = -\frac{1}{N_s} \sum_{i=1}^{N_s} X_s^i \cdot \log(p(X_s^i)) + (1 - X_s^i) \cdot \log(1 - p(X_s^i))$$

Here, $p(X_s^i) = E_s(D^s(X_s^i))$ is the source reconstruction probability of X_s^i by the source auto-encoder with N_s number of source sample. On the other hand, we design a target auto-encoder loss function considering the weighted sum of self reconstruction loss and the class-wise MMD loss as follows.

$$\mathcal{L}_t = \mathcal{L}_r + \beta \cdot \mathcal{L}_{cws_MMD}$$

$$\mathcal{L}_r = -\frac{1}{N_t} \sum_{i=1}^{N_t} X_t^i \cdot \log p(X_t^i) + (1 - X_t^i) \cdot \log(1 - p(X_t^i))$$

Here β is a scaling parameter that determines the relative importance between the class-wise MMD loss (\mathcal{L}_{cws_MMD}) and the reconstruction loss (\mathcal{L}_r).

2) *Class-wise MMD (cws-MMD) loss for Simultaneous Learning*: The MMD [29] loss function calculates the statistical distance between the centroids of the source and the target feature distributions. So, The MMD loss function can be quantified as,

$$\mathcal{L}_{MMD}(\hat{X}_s, \hat{X}_t) = \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} \hat{X}_s^i - \frac{1}{n_t} \sum_{i=1}^{n_t} \hat{X}_t^i \right\|^2$$

Where \hat{X}_s^i and \hat{X}_t^i are the i^{th} feature representation for source and target data respectively.

In designing our cws-MMD loss, we intend to calculate the divergence between class conditional probability distribution $\mathcal{P}(X|Y)$ between the source and the target datasets. As a result, the divergence of class conditional probability can be approximated by calculating the distance between the centroids of corresponding source and the target class. Then, cws-MMD loss function can be defined as,

$$\mathcal{L}_{cws_MMD}(\hat{X}^s, \hat{X}^t) = \frac{1}{C} \sum_{k=1}^C \left\| \frac{1}{n_s^k} \sum_{i=1}^{n_s^k} \hat{X}_{k,i}^s - \frac{1}{n_t^k} \sum_{i=1}^{n_t^k} \hat{X}_{k,i}^t \right\|^2 \quad (2)$$

Here, C is the maximum number of classes in the source and the target data. This new class-wise loss function aligns the source and target feature spaces according to the classes rather than the whole dataset blindly. This ensures the feature space alignment class by class for between the source and target domains.

3) *Training and fine-tuning the Classifier*: With the trained source and the target auto-encoders, we have a fine-tuned domain invariant feature space that can be used for semi-supervised (very few target labels) transfer learning. First, we train the classifier network with already trained feature space \hat{X}_s from the source encoder as features and the labels of the source domain samples as the target. We then freeze the learned classifier network and use it to predict the classes of the target feature representation \hat{X}_t in a semi-supervised fashion. The objective function of learning for the classifier network is,

$$\min_{f_c} \mathcal{L}_c[Y^s, f_c(\hat{X}_s)]$$

Here, the classifier network is represented by $f_c(\cdot)$. We use categorical cross-entropy loss for optimization of classifier network. The complete training procedure of our proposed SAEDA framework is presented in Algorithm 1.

C. Why SAEDA Works

Considering our auto-encoder based framework depicted in Eq. (1),(2), and Fig. 2, we can have the following assumptions:

a) *Hypothesis 1*: : The bottleneck feature distribution of the same classes in the source and the target domains are statistically similar. Formally, if $C_s = C_t$ then, $E_s(X_{C_s}^s) \cong E_t(X_{C_t}^t)$ where C_s and C_t are the class label of source and target domain data respectively.

b) *Hypothesis 2*: : The bottleneck feature distribution of different classes are statistically dissimilar. If $C_s \neq C_t$ then, $E_s(X_{C_s}^s) \not\cong E_t(X_{C_t}^t)$.

Algorithm 1: Training Procedure of Auto-encoder based Domain Adaptation (SAEDA)

Input : Labeled Source Domain, $\mathcal{D}^s = \{X^s, Y^s\}$,
Labeled Target Domain, $\mathcal{D}_l^t = \{X_l^t, Y_l^t\}$,
Unlabeled Target Domain, $\mathcal{D}_u^t = \{X_u^t\}$, model
parameter β , number of classifier layers c_l , and
bottleneck space size b

Output: Prediction class labels of unlabeled target domain,
 $\mathcal{D}_u^t = \{X_u^t\}$

- 1 Match the number of samples by class in \mathcal{D}^s and \mathcal{D}_l^t by randomly resampling the smaller class-domain;
 - 2 Sort \mathcal{D}^s and \mathcal{D}_l^t by class. Initialize the source and the target auto-encoder weights randomly;
 - 3 Set the loss function \mathcal{L}_s and \mathcal{L}_t to the source and the target auto-encoder;
 - 4 **repeat**
 - 5 | Train source and target auto-encoders with $\{X^s, X^s\}$ and $\{X_l^t, X_l^t\}$ respectively.
 - 6 **until** \mathcal{L}_s and \mathcal{L}_t converges;
 - 7 Take only Encoder part of source AE network and append Classifier network with it;
 - 8 Freeze the Encoder and randomly initialize the Classifier;
 - 9 **repeat**
 - 10 | Train Encoder + Classifier network with X^s, Y^s ;
 - 11 **until** Test Loss converge;
 - 12 Take target encoder and cascade with classifier network;
 - 13 Freeze target encoder part of the network;
 - 14 **repeat**
 - 15 | Train target encoder + classifier network with labelled target data, $\{X_l^t, Y_l^t\}$;
 - 16 **until** test loss converge;
 - 17 Predict the label of the target unlabelled target data, $\mathcal{D}_u^t = \{X_u^t\}$ with encoder + classifier network;
-

c) Hypothesis 3: Let denote X_C is the set of feature samples of both source and target domain of class C . For each C , there exist a domain encoder $E^*(X_C)$ and a decoder $D^*(E^*(X_C))$, and

$$\lim_{X_C \rightarrow \infty} \frac{1}{X_C} \mathcal{L}_{cws.MMD}(P_{\hat{X}_s \rightarrow t}(\cdot|C, V), P_{X_C}(\cdot|C, V)) = 0 \quad (3)$$

Where $P_{\hat{X}_s \rightarrow t}(\cdot|C, V)$ is the conditional probability of feature embedding transfer from source to target with respect to the class C and the domain invariant feature space V . Star (*) notation indicates to include both of the source and target domains.

If the bottleneck feature space dimension is set critically, the optimization loss function in Equations 3 and 4 will satisfy the ideal domain adaptation property in equation 9. This will match the class-wise domain distribution between the source and the target domain. In this case, it will only hold the domain invariant feature space representation V for every class and discard the domain-specific feature space representations U^s and U^t . This extraction of domain invariant information V in the bottleneck layers enables efficient distribution matching between source and target domain in the corresponding classes. In Fig 4-I, the decision boundary (the solid red line) is only trained on the source data which easily violates the feature space of target data.

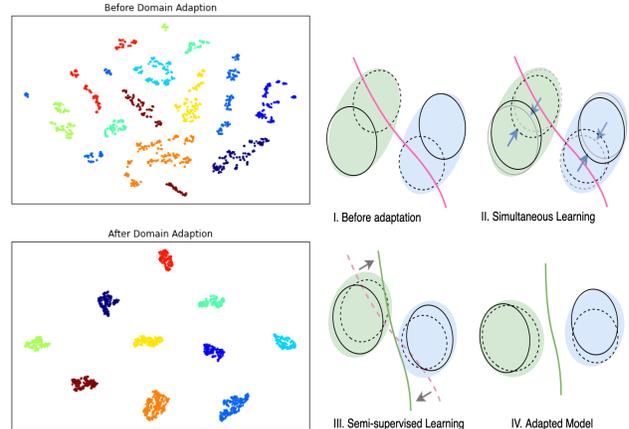


Fig. 3: The t-SNE visualization of the target domain before domain adaptation (top) and after domain adaptation (bottom). Different colors represent different classes of floor surface classification problem.

TABLE I: Vayyar Imaging Radar and Hypersen Solid State Lidar Comparisons.

Parameters	Radar	Lidar
Cost	\$599 USD	\$699 USD
Resolution	450×60	320×240
FPS	8	17
Power Supply	4.5V	9-12V
Current	0.4-0.9 A	0.25 - 1.2 A

With a simultaneous learning framework, the source and the target feature spaces align one with another (Fig 4-II). Semi-supervised fine-tuning with small numbers of target labeled data enables decision boundary to align according to the newly aligned source and the target feature spaces (Fig 4-III). Eventually, our SAEDA model generates a robust decision boundary (Fig 4-IV).

III. EXPERIMENTAL SETUP

Our sensor suite consists of an Imaging Radar (Vayyar Imaging Radar Walabot 60 GHz) and a Solid-State Lidar (Hypersen Solid-State Lidar). We align both of the sensors together focusing on a single region of interest (ROI) as show in Fig. 1. Table I provides detailed comparisons between Lidar and Radar, that prove that Radar is much efficient in terms of cost, computation (lower resolution and Frame per seconds (FPS)), and power consumption. We integrate our prototype Lidar-Radar pair with Vacuum Cleaning Robot and Mini Cheetah via USB for both power supply and data streaming.

We collected two different datasets in the lab environment and used an existing dataset which is described below:

TABLE II: Relation of sandpaper grit and surface roughness.

Sandpaper grit	120	240	320	500	1000
Roughness (μm)	59.5	30.0	23.1	15.1	9.2

- 1) **Floor Dataset:** We integrated our Lidar and Radar sensor suite in an automatic vacuum cleaning robot (as shown in Fig. 1) and collected simultaneous data on different floor types consists of (1) wooden floor, (2) tiles, (3) thin carpet, (4) thick carpet, (5) leather, (6) wet wooden floor, (7) wet tiles and (8) wet leather. The distance from the ground to the sensor is 36 cm. The data collection has been done for 5 minutes for each of the surface each, in a total, 40 minutes of collected data. We consider a 5 seconds window with 40% overlap for classification.
- 2) **Sandpaper Dataset:** We used sandpaper with various grits as the surface roughness according to the ISO 6344 international standard. The used types of sandpaper used and their respective surface roughness as presented in Table II. To measure the surface roughness, a PMMA support was fixed firmly at a distance of 25 cm from the sensor surfaces, and sandpapers of different roughness were fixed to it. We performed the data collection 30 times by change angels, in total 30,000+ frames for each sensor (Lidar and Radar).
- 3) **LAMAR Human Activity Dataset:** This is an existing dataset consists of three imaging modalities (video camera, millimeter-wave Radar and Lidar) and 6 student volunteers (graduate, undergraduate and high school students) with 7 different activities ("bending", "check_watch", "call", "single_wave", "walking", "two_wave", "normal_standing") in 3 different rooms. [30], [31]

A. Implementation Details

We implement *SAEDA* with TensorFlow and Keras. We keep the data pre-processing and the batch size consistent all over the data-sets. We keep the architecture of the auto-encoders symmetric along with the bottleneck layer as shown in Fig 2(b). As a result, the topology of decoder layers is a reverse of the encoder layers. We use upsampling 2×2 layers in the decoder modules instead of maxpool 2×2 layer in the encoder modules to keep the symmetry. There are three convolution layers and one fully connected layer in every encoder and the decoder module. We use the convolution layer as the first layer of the encoder and the second to the last layer of the decoder having 16 filters and a filter size of 2×2 . In case of the second layer of the encoder and the third to the last layer of the decoder, we use another convolution layer having 32 numbers of filters and a filter size of 3×3 . The bottleneck layer consists of 100 neurons with relu activation. Finally, we use an output convolution layer as the last layer of the decoder module of filter size 3×3 and the number of filters same as the number of channels in the images. We also use a dummy

TABLE III: Surface Roughness Estimation Result comparisons on Sandpaper dataset. Surface roughness has been estimated as (Microns) μm . R^2 score has been used to calculate the goodness of the estimation performance

Method	Radar	Lidar	Radar →Lidar	Lidar→ Radar	Radar + Lidar
CORAL [21]	0.8350	0.8532	0.9021	0.9143	0.9235
VADA [34]	0.8135	0.8298	0.9147	0.9276	0.9421
CGAN [18]	0.7064	0.7673	0.7661	0.9024	0.9143
APE [35]	0.9153	0.9312	0.9624	0.9701	0.9694
SAEDA	0.9848	0.8995	0.9976	0.9972	0.9986

layer between the bottleneck layers to maintain the source and target auto-encoders in a single graph so that we can train both simultaneously (as shown in Fig 2(b)). The dummy layer consists of a unity weight with zero activation function so that it does not affect any functionality of our system. We use Adam [32] optimization function with learning rate of 1×10^{-4} to optimize the network. We optimize the learning rate and the parameter β using hyper-parameter tuning. The final value of β is 0.25. We run our *SAEDA* model on a server having Nvidia GTX GeForce Titan X GPU and Intel Xeon CPU (2.00GHz) processor with 12 Gigabytes of RAM.

IV. RESULTS

We implement **four baseline methods:** Deep Correlation Alignment (CORAL) [33], Virtual Adversarial Domain Adaptation (VADA) [34], Conditional GAN (CGAN) [13] and Attract, Perturb, and Explore (APE) [35].

A. Classification Results: Activity and Floor Recognition

Fig.6(a) and Fig. 6(b) show performance comparisons of *SAEDA* framework with baseline methods. It can be depicted that *SAEDA* outperforms baseline methods significantly providing on average 4.5% and 3.2% improvements of accuracies than the nearest baseline methods for activity recognition and floor surface recognition respectively. Fig. 5 shows the confusion matrix of activity recognition and floor surface recognition for both Lidar → Radar and Radar → Lidar domain adaptations. Fig. 5 shows that Radar → Lidar provides extremely low accuracy comparing to Lidar → Radar adaptations. The accuracy differences are expected as we know 60 GHz Radar is highly absorbed by water providing significant changes in signals while Lidar signals are partially absorbed by water. On the other hand, Radar → Lidar provides better accuracy than Lidar → Radar due to Lidar's higher accuracy in position measurement (precision level is 1 cm for Hypersen Solid-State Lidar) than Radar.

B. Regression Results: Surface Roughness Characterization

We substitute softmax layer of *SAEDA* framework to a linear regression layer to create a regression problem for surface roughness estimation. Table III shows the comparisons of R^2 (*R - Squared*) and Mean Squared Error (MSE) metrics of the different algorithms in estimating surface roughness. *R-squared* quantifies the relative percentage measure of the

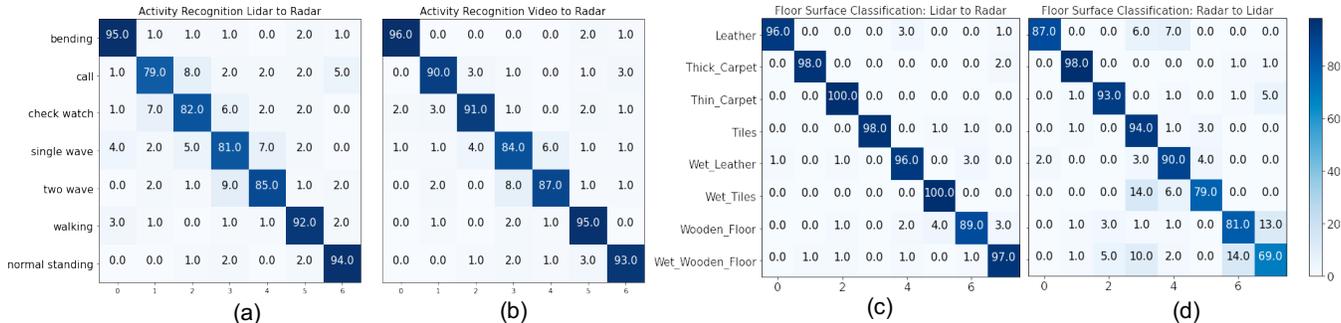


Fig. 5: Confusion Matrix of activity recognition accuracy for (a) Lidar \rightarrow Radar and (b) Video \rightarrow Radar domain adaptation on LAMAR dataset as well as floor surface classification for (c) Lidar \rightarrow Radar and (d) Radar \rightarrow Lidar on our dataset

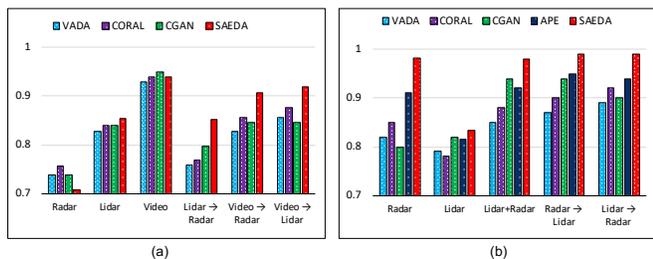


Fig. 6: (a) Activity Recognition Result Comparisons with Baseline (b) Floor surface recognition result Comparisons with Baseline methods

variance of dependent variables from the model, which can be defined as $R^2 = 1 - \frac{\text{sum of squares of residuals}}{\text{total sum of squares}}$. It can be easily viewed that the Radar sensor outperforms Lidar in detecting surface roughness. Also, *SAEDA* framework improves the domain adaptation significantly than the baseline methods.

C. Outdoor Test using Mini-Cheetah Robot

To see the feasibility of our roughness estimator in the four-legged robot application, we integrate our Radar and Lidar sensor suite with Mini-Cheetah and run an outdoor experiment. In this test, Mini-Cheetah moves around different surfaces which include thick grass, leather carpet, leaf, and thin grass (Fig. 7). The locomotion controller of [36] is used and there is no feedback from the vision sensors in this test. We intentionally decouple the walking control from the estimation to watch that the terrain roughness is reasonably estimated even under the notable disturbance coming from the general walking behavior. The test results are summarized in Table IV. It is visible that leaves are the softest surface in the wild while thick grass can be considered as the roughest surface as per our feasibility study.

V. DISCUSSION

Figure 3 shows the t-SNE [37] feature representation of target domain features before and after domain adaptation. Different colors represent different classes in the target domain. From the t-SNE plot, it is evident that feature representation space is not well clustered before domain adaptation which leads to misclassification and hence lower

TABLE IV: Mini-Cheetah Robot Outdoor Test Result.

Surface	Grass	ThickCarpet	LeatherCarpet	Leaves
Roughness (μm)	30.5 ± 3.5	19.5 ± 2.3	15.2 ± 2.6	10.6 ± 1.4



Fig. 7: Mini-Cheetah walks over different surfaces (grass, carpet, leaf surfaces) while carrying our sensor suite.

VI. CONCLUSION

SAEDA is capable of learning domain invariant feature space between simultaneously collected noisy and error-free sensor signals, enabling the development of the semi-supervised transfer learning framework. Experimental evaluation on three distinct recognition tasks using collected and public datasets, and feasibility study provides ample evidence that *SAEDA* is capable of amplifying any robotic sensing performance with the involvement of least possible hardware and computational complexity. This framework, revolutionizing the idea of hardware integration virtualization, is able to provide existing deployed sensors-integrated autonomous system such capabilities that only could be achieved using expensive hardware modifications and computationally complex infrastructure or replacement of the latest version of mobile robotic systems.

REFERENCES

- [1] J. a. Guerreiro, D. Sato, S. Asakawa, H. Dong, K. M. Kitani, and C. Asakawa, "Cabot: Designing and evaluating an autonomous navigation robot for blind people," in *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*, ser. ASSETS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 68–82.

- [2] T. Yonezawa, H. Yamazoe, A. Utsumi, and S. Abe, "Gazerboard: Gaze-communicative guide system in daily life on stuffed-toy robot with interactive display board," in *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems, September 22-26, 2008, Acropolis Convention Center, Nice, France*. IEEE, 2008, pp. 1204–1209. [Online]. Available: <https://doi.org/10.1109/IROS.2008.4650692>
- [3] Y. Zhang and Y. Meng, "A decentralized multi-robot system for intruder detection in security defense," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010, pp. 5563–5568.
- [4] M. A. Arain, V. H. Bennetts, E. Schaffernicht, and A. J. Lilienthal, "Sniffing out fugitive methane emissions: autonomous remote gas inspection with a mobile robot," *The International Journal of Robotics Research*, vol. 0, no. 0, p. 0278364920954907, 2020.
- [5] C. Debeunne and D. Vivet, "A review of visual-lidar fusion based simultaneous localization and mapping," *Sensors*, vol. 20, no. 7, 2020. [Online]. Available: <https://www.mdpi.com/1424-8220/20/7/2068>
- [6] H. A. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," *CVPR*, pp. 12 697–12 705, 2019.
- [7] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 4490–4499.
- [8] B. Yang, W. Luo, and R. Urtaasun, "PIXOR: real-time 3d object detection from point clouds," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 7652–7660.
- [9] R. Weston, S. Cen, P. Newman, and I. Posner, "Probably unknown: Deep inverse sensor modelling radar," in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5446–5452.
- [10] P. Kaul, D. De Martini, M. Gadd, and P. Newman, "Rss-net: Weakly-supervised multi-class semantic segmentation with fmcw radar," *arXiv preprint arXiv:2004.03451*, 2020.
- [11] G. Kim, Y. S. Park, Y. Cho, J. Jeong, and A. Kim, "Mulran: Multimodal range dataset for urban place recognition," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 6246–6253.
- [12] G. Kim and A. Kim, "Scan context: Egocentric spatial descriptor for place recognition within 3d point cloud map," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 4802–4809.
- [13] Y. S. Park, J. Kim, and A. Kim, "Radar localization and mapping for indoor disaster environments via multi-modal registration to prior lidar map," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2019, Macau, SAR, China, November 3-8, 2019*. IEEE, 2019, pp. 1307–1314.
- [14] H. Yin, Y. Wang, L. Tang, and R. Xiong, "Radar-on-lidar: metric radar localization on prior lidar maps," 2020.
- [15] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7167–7176.
- [16] G. Csurka, "Domain adaptation for visual applications: A comprehensive survey," *arXiv:1702.05374, arXiv preprint*, 2017.
- [17] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko, "Simultaneous deep transfer across domains and tasks," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [18] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [19] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *European conference on computer vision*. Springer, 2016, pp. 443–450.
- [20] W. Zellinger, T. Grubinger, E. Lughofer, T. Natschläger, and S. Saminger-Platz, "Central moment discrepancy (cmd) for domain-invariant representation learning," *arXiv preprint arXiv:1702.08811*, 2017.
- [21] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," *arXiv preprint arXiv:1502.02791*, 2015.
- [22] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Unsupervised domain adaptation with residual transfer networks," in *NIPS*, 2016, pp. 136–144.
- [23] Y. Yao, Y. Zhang, X. Li, and Y. Ye, "Heterogeneous domain adaptation via soft transfer network," in *Proceedings of the 27th ACM International Conference on Multimedia*, 2019, pp. 1578–1586.
- [24] J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. Wortman, "Learning bounds for domain adaptation," in *NIPS*, 2008, pp. 129–136.
- [25] B. Katz, J. D. Carlo, and S. Kim, "Mini cheetah: A platform for pushing the limits of dynamic quadruped control," 2019, pp. 6295–6301.
- [26] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2208–2217.
- [27] B. Sun, J. Feng, and K. Saenko, "Return of frustratingly easy domain adaptation," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [28] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [29] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola, "A kernel method for the two-sample-problem," *Advances in neural information processing systems*, vol. 19, pp. 513–520, 2006.
- [30] M. Alam, M. M. Rahman, and J. Widberg, "Palmar: Towards adaptive multi-inhabitant activity recognition in point-cloud technology," *IEEE International Conference on Computer Communications*, 2021.
- [31] C. Benedek, B. Gálai, B. Nagy, and Z. Jankó, "Lidar-based gait analysis and activity recognition in a 4d surveillance system," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 1, pp. 101–113, 2018.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [33] B. Sun and K. Saenko, "Deep CORAL: correlation alignment for deep domain adaptation," *CoRR*, vol. abs/1607.01719, 2016. [Online]. Available: <http://arxiv.org/abs/1607.01719>
- [34] R. Shu, H. H. Bui, H. Narui, and S. Ermon, "A dirt-t approach to unsupervised domain adaptation," 2018.
- [35] T. Kim and C. Kim, "Attract, perturb, and explore: Learning a feature alignment network for semi-supervised domain adaptation," in *European Conference on Computer Vision*. Springer, 2020, pp. 591–607.
- [36] D. Kim, J. Di Carlo, B. Katz, G. Bledt, and S. Kim, "Highly dynamic quadruped locomotion via whole-body impulse control and model predictive control," *arXiv preprint arXiv:1909.06586*, 2019.
- [37] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.