

RGB-X Classification for Electronics Sorting

FNU Abhimanyu¹, Tejas Zodage¹, Umesh Thillaivasan², Xinyue Lai¹, Rahul Chakwate¹,
Javier Santillan², Emma Oti², Ming Zhao², Ralph Boirum¹ Howie Choset¹, Matthew Travers¹

Abstract—Effectively disassembling and recovering materials from waste electrical and electronic equipment (WEEE) is a critical step in moving global supply chains from carbon-intensive, mined materials to recycled and renewable ones. Conventional recycling processes rely on shredding and sorting waste streams, but for WEEE, which is comprised of numerous dissimilar materials, we explore targeted disassembly of numerous objects for improved material recovery. Many WEEE objects share many key features and therefore can look quite similar, but their material composition and internal component layout can vary, and thus it is critical to have an accurate classifier for subsequent disassembly steps for accurate material separation and recovery. This work introduces RGB-X, a multi-modal image classification approach, that utilizes key features from external RGB images with those generated from X-ray images to accurately classify electronic objects. More specifically, this work develops Iterative Class Activation Mapping (iCAM), a novel network architecture that explicitly focuses on the finer-details in the multi-modal feature maps that are needed for accurate electronic object classification. In order to train a classifier, electronic objects lack large and well annotated X-ray datasets due to expense and need of expert guidance. To overcome this issue, we present a novel way of creating a synthetic dataset using domain randomization applied to the X-ray domain. The combined RGB-X approach gives us an accuracy of 98.6% on 10 generations of modern smartphones, which is greater than their individual accuracies of 89.1% (RGB) and 97.9% (X-ray) independently. We provide experimental results³ to corroborate our results.

I. INTRODUCTION

Daisy [1], [2] and Dave [3], [4], Apple’s existing disassembly robots, disassemble electronic devices and components to enable the recovery of precious materials like rare earth elements, steel, and tungsten. Diverse object disassembly involves hundreds of decisions throughout the process, and object classification is the first major decision in the pipeline. In practice, precisely classifying electronics objects is challenging as different objects of similar product line often have subtle visible differences leading to incorrect classifications that can cause sub-optimal material recovery if material compositions differ between the actual and predicted class. Conventional machine vision has to be highly accurate for recycling, however to classify all desired WEEE objects for disassembly using one line also requires that the system be capable of handling high variations seen in each class of multi-class WEEE such as cracks, bends, missing or replaced components, or other deformations that are difficult to predict

¹Robotics Institute, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, USA abhiman2@andrew.cmu.edu, choset@andrew.cmu.edu, mtravers@andrew.cmu.edu

²Apple Inc.

³Experimental work done at Biorobotics Lab, Robotics Institute, Carnegie Mellon University

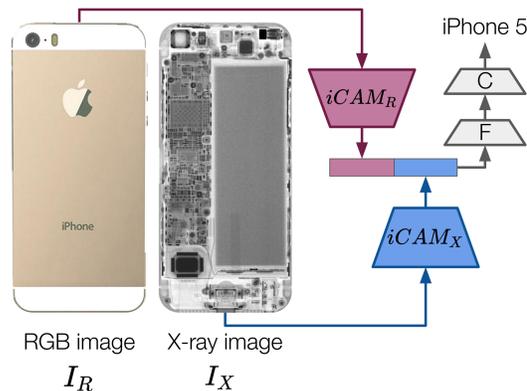


Fig. 1: Our approach fuses feature maps of RGB I_R and X-ray image I_X in order to capture both internal and external features of electronics devices. We develop a novel architecture, $iCAM$ for suitable feature maps extraction of electronics images and show that the RGB-X fusion along with $iCAM$ gives classification accuracy as high as 98.8%.

and that inhibit conventional classification methods (e.g., RGB-based methods). These variations eliminate methods like weighing when numerous classes differ by only a few grams. Thus, this paper addresses challenges in classifying WEEE as they have subtle differences on the outside forcing a classifier to focus on extremely specific details on the outside or look inside for other distinguishing features.

In this work, we present a novel attention based iterative strategy, Iterative Class Activation Mapping (iCAM), to localize and then extract key features of relevant components from inside and outside of the devices for classification. This effectively guides the models to focus on the appropriate features for better decision making.

Also, to leverage both outside and inside features of electronic objects, we propose RGB and X-ray classification modalities individually and then uniquely combine them for their classification. Though there has been previous work on multi-domain classification, to the best of our knowledge, this work is the first to combine RGB and X-ray domains to perform classification for recycling purposes. This combination helps find common key features among the internal and external parts of the objects, making the classification robust to conditions like lighting, or wear and tear of the device.

The single-mode (RGB mode or the X-ray mode), as well as the joint training phase, needs an annotated dataset for training. Unlike other areas of image classification, used electronics lack a well-annotated dataset for training an image classifier robustly, especially in the X-ray domain.



Fig. 2: Examples of the training and the testing data for the RGB and the X-ray images. From left to right: synthetic RGB images (training), synthetic X-ray images (training), real RGB images (testing), real X-ray images (testing).

The creation of annotated X-ray data is tedious and requires expert guidance. Thus, in this work, we propose a data randomization pipeline to generate synthetic X-ray data for training. We create synthetic X-ray images by projecting 3D-Computed Tomography (CT) scans into a 2D space and randomizing the relevant parameters like orientation, intensity settings, noise and background.

The paper is organized as follows: Section II introduces the prior work reported for electronic objects' image classification. Section III provides a detailed description on the dataset, iterative class activation mapping (iCAM) architecture and the training strategy. Section IV describes the experimental results on RGB-X classification. Section V discusses the conclusion and the future work.

II. PRIOR WORK

This section discusses the existing WEEE classification methods. It also discusses the prior work regarding the dataset generation process, fine-grained image classification and multi-modal classification.

A. Classification for WEEE sorting

WEEE is an existing problem where both classical techniques as well as deep learning architectures have been used to recognize objects. UNU-KEYS [5] is a classification method that is being used to classify WEEE by using attributes such as average weight, material compositions and end-of-life characteristics. This works well for coarse level classification, but does not work well for similar looking devices which have subtle feature differences. To overcome that, state of the art deep learning methods have also been used to accurately classify WEEE items. Standard network architectures like Faster Squeeze-Net [6] and YOLOv3 [7] were used to sort electronics. The Swedish company, Refind, in collaboration with the Danish Institute of Technology has also demonstrated the feasibility of sorting different types of printed circuit boards, mixed electronic scrap, and batteries by capturing color images and applying a deep learning architecture [8], [9]. Use cases include applying deep learning to sort different types of coin cells [10], and additional imaging inputs such as thermal and X-ray imaging are being explored for WEEE recognition [11], [12]. While the prevalence of deep learning has allowed for advances in classification, these methods are still primarily used for

component level classification and they use a single mode for classification, therefore making models sensitive to real world scenarios like lighting, noise, and deformations.

B. Domain Randomization

One of the obstacles of using a machine learning algorithm is the lack of readily available, large, annotated datasets. The machine learning community has developed the concept of domain randomization (DR) to synthetically generate large datasets instead of collecting and labeling data manually. DR aims to add large variability to the synthetic data generation process, attempting to capture the real-world variations.

Svetozar and Jianhong [13] show a useful analysis of important parameters to generate synthetic data for classification for RGB images. DR has also been used in non-classification tasks. Tobin *et. al.* [14] used DR to synthetically generate images of geometrical shapes for an autonomous robotic picking task. Additional applications include autonomous drone flight [15], object detection [16], viewpoint estimation [17], and human pose estimation [18]. These works demonstrate the application of DR in the RGB domain, but they do not explore their use in other domains like X-rays.

C. Fine-grained image classification

Fine-grained image classification captures discriminative features of fine-grained objects and uses those features in image classification. In this paradigm, a localization sub-network is designed for locating the key regions of interest followed by a classification network used for recognition. Previous works on this paradigm depended on adding dense part annotations for accurate part localization and then using them for image recognition [19], [20], [21]. Recent works require only image labels for accurate part localization for training the recognition sub-network [22], [23], [24]. Finally, additional techniques such as attention mechanisms [25] and multi-stage strategies [26] perform the joint training of the integrated localization-classification sub-networks.

D. Multi-domain classification

Multi-domain learning aims to improve the classification performance in general domains by making full use of the information in each domain. Prior work has focused on extracting hand-crafted features such as scale-invariant

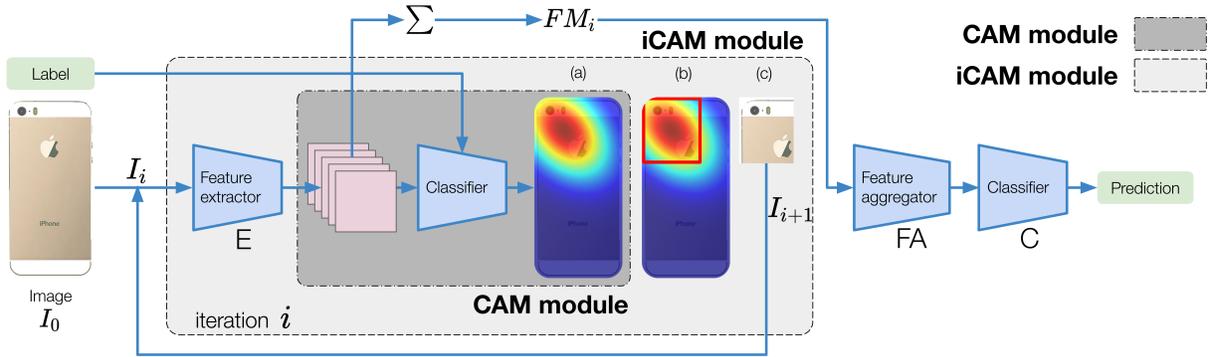


Fig. 3: Classification network architecture with SP1 as the input image I_0 . For iteration i , an image I_i is passed as an input to the $iCAM$ module. Internally, the $iCAM$ module extracts feature maps (FM_i) of image I_i and extracts the region of activation (i.e., the region most important for classification) using the CAM module. Based on the importance regions, an image I_{i+1} is cropped from I_i , which is used for the iteration $i+1$. For the final prediction, we aggregate the feature maps in the feature aggregator (FA) and pass it through a classifier block to generate a prediction. Subimages (a), (b) and (c) shows the heat map, the most important region of the image given the heat map, and the cropped image from the heat map.

feature transform (SIFT) from the multi-domain (RGB-D) image to get a combined image descriptor [27]. Blum et al. proposed an RGB-D descriptor that relies on a k-means-based feature learning approach [28]. However, the hand-crafted features are often dataset-specific and require a strong understanding of domain-specific knowledge [29]. To reduce the dependency on hand-crafted features, machine learning techniques are explored, and one of the most common techniques is using network parameter sharing [30].

III. APPROACH

This section discusses the X-ray dataset generation process, the network architecture used for classification, and the joint training strategy used for combining the X-ray mode with the existing RGB mode for a combined RGB-X classification. In both these domains, we conduct experiments on 10 classes of modern smart phones (SP) classes.

A. Datasets

The 10 classes of modern smart phones will be referred to as SP1, SP2, SP3, SP4, SP5, SP6, SP7, SP8, SP9, SP10 in this work. For the RGB domain, the test dataset consists of 799 real RGB images. RGB training set includes a collection of synthetic and real images. The training dataset consists of a combination of 1000 synthetic and 250 real images per class. We collect RGB images using a Microsoft Azure Kinect camera mounted over a conveyor belt. We load the mesh of the phones in Blender [31] and use DR to create synthetic dataset by varying color, texture, and camera intrinsic and extrinsic parameters as suggested in [14].

For the entire X-ray domain, the test dataset consists of 749 real X-ray images and around ≈ 25000 synthetic X-ray images for training. Only synthetic X-ray images are used for training to demonstrate the efficacy of domain randomization and because of the lack of enough real X-ray present for each class. Test images are collected using an X-ray system equipped with 300kV micro focus and 180kV nano focus tubes, and a cesium iodide (CsI) detector. Both tubes are

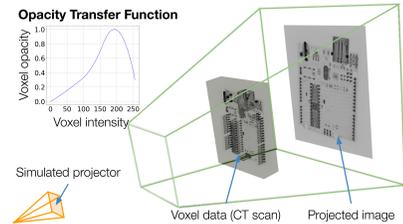


Fig. 4: X-ray Domain Randomization Setup. To generate synthetic X-ray images, we project CT scans onto a plane, simulating the process of taking an X-ray image of an electronic object. We show our setup on a simple electronic device due to privacy reasons tied to the phone image.

used to generate 2D and 3D images with resolutions ranging from $57.5 \mu\text{m}$ to $100 \mu\text{m}$. The 2D and 3D images taken with the 300 kV tube use a peak energy of 220 kV with a current ranging from 200 to 300 μA , while images acquired with the 180 kV tube use 160 kV and 280 μA . For training, we utilize DR in order to generate a large scaled labelled dataset. To generate X-ray training data, we initially generate 3D voxel models of 2-3 devices per class using off the shelf CT machines. CT scans consist of a scalar intensity assigned to each voxel based on the structure of the device. The voxel intensities are mapped to an opacity value and projected onto a plane using a simulated projector model to render X-ray images. The voxel intensity to voxel opacity mapping is known as a transfer function. A variety of X-ray machine settings can be simulated by the shape of the transfer function as shown in Fig 4. We use VTK library [32] to randomize the transfer function, and the intrinsic and extrinsic parameters of simulated projectors. Randomizing the transfer function arbitrarily can undesirably produce completely saturated images by assigning the same value to each pixel. To avoid this, we seed the randomization with 5 manually selected transfer function parameters. These manual seeds were selected using the visualization tool 3D Slicer [33]. We add randomness to these parameters and only select rendered images that have more than 100 SuperPoint [34] keypoints. Sample training

and testing images from the RGB and X-ray dataset are shown in Fig 2.

B. Network architecture and training strategy

This section describes the fine-grained classification network used for training the RGB and the X-ray mode, and the training strategy used in the single training and joint training phases. Let RGB image (I_R) or X-ray image (I_X) be of dimension $H \times W \times C$ in the labeled dataset available for training the classification model, and y_i be the image label in one-hot encoding - *i.e.*, $y_i \in R^m$ is a vector of dimensionality m (the number of classes) with $y_i^{j=l} = 1$ and $y_i^{j \neq l} = 0$ for the position l denoting the image label. We train the model using a three-stage approach: first training the 2 stream networks individually, followed by a joint fine-tuning stage.

1) *Iterative class activation mapping (iCAM)*: As shown in Fig 3, our network has 4 modules: feature extractor module (E), class activation mapping module (CAM), feature aggregator module (FA), and classifier module (C). The feature extractor module (E) and the CAM module are stacked together in an iterative manner followed by the feature aggregator module and classifier module.

In our method, we use class activation mapping (CAM) [35] using the global average pooling (GAP) to indicate the discriminative image regions used by CNNs to identify the respective classes. For a given image $I \in \mathbb{R}^{H \times W \times C}$, let $FM(x, y) = \sum_{j=1}^k fm_j(x, y)$ represent the activation of the last convolutional layer at spatial location (x, y) with k channels. θ_C^c denotes the weight corresponding to class c in the classifier module C . We define m_c as the class activation map for class c as:

$$\begin{aligned} fm_{0\dots k}(x, y) &= f_{\theta_E}(I), \\ m_c(x, y) &= \sum_{j=1}^k \theta_C^c fm_j(x, y). \end{aligned} \quad (1)$$

Thus, $m_c(x, y)$ is an unnormalized probability distribution directly indicating the importance of the activation at spatial location (x, y) . More details on CAM can be found in [35].

In our experiments, InceptionNet-v3 [36] up to the Mixed7c layer is chosen as the feature extraction module with θ_E as a learnable parameter. The feature map (FM_i) and the class activation map ($m_{c,i}$), for every iteration i , is calculated using Eq.1. The activation map is upsampled to the original image size $H \times W \times C$, and then a square region around the most dominant pixel of the map is selected resulting in image I_{i+1} . The cropped image is then upsampled to the original image dimension $H \times W \times C$ using bilinear upsampling function in Pytorch. After running the feature extraction and the CAM module for a pre-defined number of iterations n , $F_{0..n}$ is passed to the feature aggregator module followed by the classifier module to get the final classification score S :

$$\begin{aligned} FM(x, y) &= \sum_{i=1}^n \theta_{FA} FM_i(x, y), \\ S &= \text{softmax}(f_{\theta_C}(FM)), \end{aligned} \quad (2)$$

where FM is the aggregated feature map, and θ_{FA} and θ_C are the learnable parameters of the feature aggregator and the classifier module. For our experiments, we choose $n=3$.

2) *Training strategy*: We proceed by training the iCAM network to minimize the negative log likelihood, \mathcal{L} , of the training data. While training we solve for:

$$\begin{aligned} \theta_E^*, \theta_{FA}^*, \theta_C^* &= \\ \text{argmin}_{\theta_E, \theta_{FA}, \theta_C} &\sum_{i=1}^N \mathcal{L}(f_{iCAM}(I, y_i; \theta_E, \theta_{FA}, \theta_C), y_i), \end{aligned} \quad (3)$$

We individually train the $iCAM_R$ and $iCAM_X$ for I_R and I_X images on RGB and X-ray datasets. For both the streams, the θ_E is initialized with InceptionNet-v3 weight pretrained on ImageNet and were trained for 40 epochs each. The implementation is highly parallelized and performs full-batch gradient descent using the Stochastic Gradient Descent [37] optimizer in the Pytorch Autograd library [38], with a batch size of 16 with a learning rate of 0.001. At the end of this individual training stage, $iCAM_R$ and $iCAM_X$ have different sets of weights which are later fused in the joint training stage.

Once both domains are trained individually, we use the parameter sharing method to combine the information of both the domains as shown in Fig1. During the joint-training, the softmax activation is discarded and the output from the classification layers are concatenated. Their individual responses FM_R from $iCAM_R$, and FM_X from $iCAM_X$, are fused and fed through an additional stream $FS([g_R, g_X])$ with parameters θ_{FS} . This fusion network again ends in a softmax classification layer. During the joint-training phase, $\theta_{E,R}$, $\theta_{FA,R}$, $\theta_{E,X}$ and $\theta_{FA,X}$ are frozen and only $\theta_{C,R}$, $\theta_{C,X}$ and θ_F are open to training. In this stage of training, we solve:

$$\begin{aligned} \theta_{C,R}^*, \theta_{C,X}^*, \theta_{FS}^* &= \\ \text{argmin}_{\theta_{C,R}, \theta_{C,X}, \theta_{FS}} &\sum_{i=1}^N \mathcal{L}(f_{FS}([g_R(I_R), g_X(I_X)], y_i). \end{aligned} \quad (4)$$

During the joint-training phase, only the weights of the fusion layer and the weights of individual classification layers are optimized, keeping all other weights from individual mode intact.

IV. ANALYSIS AND RESULTS

A. Domain randomization

In this set of experiments, we evaluate the effect of domain randomization parameters on test accuracy. Our current implementation randomizes image white noise, background color, transfer function parameters, and rendering projector position and orientation. As shown in Table I in the first experiment, we introduce noise of increasing variance in the rendered image. We observe that the best accuracy of 98.1% is obtained when $\sigma = 1000$ with varying background, transfer function, and projector position and orientation. Also, we see a drop in the accuracy with the drop in the noise variance. The presence of noise improves convergence

TABLE I: Effect of DR parameters on accuracy. Current DR has non-fixed background, non-fixed transfer function, non-fixed simulated projector position and orientation and $\sigma = 1000$ as the variance for the white noise.

Parameters	Status	Accuracy
Current DR	Sec. III	0.981
Noise	$\sigma=10$	0.934
	$\sigma=100$	0.923
	$\sigma=500$	0.957
Background	Fixed	0.922
Transfer function	Fixed	0.655
Simulated projector position and orientation	Fixed	0.315

and makes training less susceptible to local minima as shown in [39]. We also observe that fixing the projector position and orientation affects accuracy the most, lowering it down to only 31.5%. This result emphasizes the use of DR as this projector position simulates different views of a 3D object, which can not be simulated by simple 2D data augmentation. Another important factor is the transfer function, which as discussed in Sec. III, simulates different X-ray machines and X-ray machine settings. Varying the transfer function, generalizes our data over different machine distributions. Unlike the RGB domain, we avoid pattern randomization since such variations do not exist in our test data.

B. iCAM

iCAM is a fine grained classification method. For the sake of result section, we show the performance of iCAM on the modern smart phones (SP) classes mentioned in Section III and also CUB-200 [40] dataset. For comparisons, we evaluate the results with baseline classification networks like VGG-19 [41], ResNet-101 [42], ResNeXt-101 [43] and InceptionNet-v3 [44]. We also compare it with the SqueezeExcitation [45] with ResNet101 as the base module. All these experiments are performed on the X-ray dataset. We use average accuracy, and inference time as the metrics to compare the results for different networks as suggested in [46]. We also compare average precision and average recall of all these network architecture to our method as suggested in [47]. These results are presented in Table II. Additionally, we compare the iCAM+InceptionNet-v3 to these baseline classification networks on the CUB-200 dataset [40], a standard fine-grained classification dataset used for bench-marking image classification algorithms.

Table II shows that the iCAM block when combined with the InceptionNet-v3 base module outperforms the baseline image classification networks. The high accuracy value is attributed to the hierarchical approach iCAM takes in order to extract and combine key features from an image. Intuitively this approach is similar to looking and extracting key features at different levels of magnification and then using all the features for classifying that image. For further analysis we calculate the heat-map using CAM as seen in Fig 6, and use them to visualize the CNN’s area of interest for similar looking classes (e.g., SP2 vs SP3). It is seen that iCAM block indeed helps in localizing key features in an X-ray image. Additionally, Table II presents that iCAM+InceptionNet-v3

TABLE II: Comparison of baseline classification networks to the iCAM network on the X-ray dataset. Learning rate: 0.001, Batch Size: 16, Weight decay: 0.1 every 7 epochs, Epochs: 40. The bolded row highlights the network architecture with the best average accuracy.

Network Architecture	Average accuracy	Precision	Recall	Infer time (ms)
InceptionNet-v3	0.933	0.907	0.935	12.3
ResNet-101	0.953	0.899	0.956	11.1
VGG-19	0.956	0.898	0.961	7.53
ResNeXt-101 (64x4d)	0.929	0.869	0.948	20.15
SE-ResNet-101	0.891	0.841	0.926	14.01
iCAM+ InceptionNet-v3	0.979	0.942	0.969	30.86

TABLE III: Comparison of baseline classification networks to the iCAM network on the CUB-200 [40] dataset. Same parameters as Table 2. The bolded row highlights the network architecture with the best average accuracy for InceptionNet-v3 and ResNet-101 feature extractor.

Network Architecture	Accuracy	Precision	Recall
InceptionNet-v3	0.692	0.697	0.693
iCAM+ InceptionNet-v3	0.797	0.791	0.801
ResNet-101	0.814	0.818	0.815
ResNeXt-101 (64x4d)	0.829	0.835	0.830
SE-ResNet-101	0.818	0.826	0.820
iCAM+ ResNet-101	0.834	0.831	0.838
VGG-19	0.778	0.781	0.779

has an inference time of ≈ 30 fps, which allows the model to run in real-time with a 30fps industrial camera. Table III, shows that when iCAM is combined with InceptionNet-v3, the network outperforms InceptionNet-v3 by 10.5%. We also combine the iCAM block with the ResNet-101 feature extractor, and the result outperforms ResNet-101 [42], and other networks like ResNeXt-101 [43], SE-ResNet-101 [45], and VGG-19 [41]. This shows the efficacy of iCAM block with other extraction modules.

We also conduct a study on the number of iterations the iCAM module is run during the training time. Table IV shows the accuracy, the time taken to train 40 epochs, and the inference time for $n = 2, 3, 5$. As shown in the table, $n = 3$ shows a considerable increase in accuracy from using $n = 2$ for far less training and inference time that $n = 5$.

A comparative study is done to select the best feature aggregator strategy where 3 different feature aggregator strategies are tested. The most common strategies used are mean, sum, and weighted mean of the available feature maps. All of these experiments are done for $n = 3$. We achieve an accuracy of 0.941 for the mean, 0.768 for the sum, and 0.979 for the weighted mean version of the feature aggregator. Based on this study, we conclude that weighted mean is the best feature aggregating strategy because it weighs different



Fig. 5: Examples of training and testing images from the CUB-200 [40] dataset. Small inter-class variance and large intra-class variance between species motivates us to use this dataset for benchmarking iCAM.

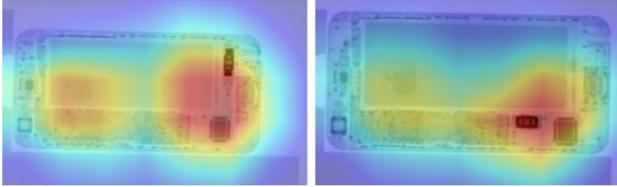


Fig. 6: Heatmap for SP2 (left) and SP3 (right) generated using CAM. SP2 has a prominent logo visible in the X-ray image, that is the dominant region for classification. SP3 has a logo and a different component, that helps identify SP3.

level of magnification differently .

C. RGB-X classification

We evaluate the multi-modal classification approach with both the single-mode iCAM+InceptionNet-v3 architectures. For this experiment, an RGB image of a particular class is paired with a random X-ray image sampled from the same class. This is done because of the lack of a streamlined RGB + X-ray sensor setup. While pairing the images, we randomly re-orient the X-ray images to prevent any bias induced due to similar orientation of the RGB and X-ray image. The total training time for RGB-X mode is 520 mins as the RGB and X-ray modes are trained in parallel. Table V shows the average accuracy, precision, recall and the inference time of each of the mode (RGB and X-ray) as well as the combined RGB-X mode. In context of a recycling plant that recycles around 200 phones per hour [2], the 0.6 percentage gain in accuracy from RGB-X ($n = 3$) over X-ray ($n = 5$) is critical as it saves 12 phones per hour. Although this adds 11ms to the inference time, the total training time is lowered to 520 mins (RGB-X) from 723 mins (X-ray). The lower training time is crucial while expanding to newer products without compromising on the recycling rate of 200 phones per hour.

As mentioned in Sec.III-B.2, for joint training we combine the deepest layers and train for them, by keeping the higher

TABLE IV: Comparison of accuracy, training time and inference time of the iCAM+InceptionNet-v3 network for different values of iteration N. Training was done for 40 epochs. The bolded row highlights the current n chosen for all the experiments.

No. of CAM iteration (N)	Average accuracy	Training time (mins)	Infer-time(ms)
n=2	0.962	325	20.11
n=3	0.979	490	30.86
n=5	0.981	723	50.96

TABLE V: Comparison of single (RGB, X-ray) mode accuracy to the multi-modal (RGB-X mode) accuracy. The bolded row highlights the mode with the best average accuracy.

Mode	Average accuracy	Average precision	Average recall	Infer time (ms)
RGB	0.891	0.932	0.863	30.86
X-ray	0.979	0.942	0.969	30.86
RGB-X	0.987	0.987	0.995	62.4

TABLE VI: Comparison of different fusion techniques. The bolded row highlights the mode with the best average accuracy.

Mode	Average accuracy	Average precision	Average recall
Slow Fusion	0.952	0.951	0.958
Fast Fusion	0.987	0.987	0.995

CNN layers frozen. This fusion approach is termed as “fast fusion” [48]. We evaluate this approach with the “slow fusion” approach also mentioned in [48]. For slow fusion, the network architecture is similar to fast fusion, except for the fact that the weights of the higher CNN layers are also trainable. This allows the higher layers to have access to more global information.

Table VI, shows the comparisons of 2 different fusion techniques, demonstrating that fast fusion performs better than the slow fusion, as fast fusion doesn’t affect the lower level activation learned by the iCAM layers.

V. CONCLUSION AND FUTURE WORK

We introduce a useful multi-modal neural network architecture for the RGB-X object for the task of WEEE classification. Our method consists of a two-stream convolutional neural network that learns to fuse information from both RGB and X-ray domains automatically before classification. We also present a novel fine-grained network architecture, which is used to train the individual streams of data. This method iteratively localizes key features in an object and uses them for classification. Our experiments explore how the various components of this algorithm achieve better classification than the baseline algorithms when tested on 10 classes of modern smart phones. For this paper, we use iPhone as our modern smart phones, as the devices and the dataset were readily available. We also present a novel domain randomization pipeline for X-ray images using CT scans, a method that reduces the burden of annotating thousands of X-ray images. As our future work, we plan to extend the work to classify smart phones from different manufacturers like Samsung, Motorola and also between inter-manufacturer models. Additionally, we plan to extend the approach to attend to multiple important regions in an activation map. We also plan to generalize the hierarchical attention approach to other interpretable methods like Grad-CAM [49].

The high accuracy as well as high inference rate of this approach make it well suited for integration into electronics sorting processes at material recovery facilities.

REFERENCES

- [1] Apple. Apple adds earth day donations to trade-in and recycling program. [Online]. Available: <https://www.apple.com/newsroom/2018/04/apple-adds-earth-day-donations-to-trade-in-and-recycling-program>
- [2] Apple. Apple expands global recycling programs. [Online]. Available: <https://www.apple.com/newsroom/2019/04/apple-expands-global-recycling-programs>
- [3] Apple. Apple commits to be 100 percent carbon neutral for its supply chain and products by 2030.
- [4] Apple. Environmental progress report. [Online]. Available: <https://www.apple.com/environment/pdf/Apple%20Environmental%20Progress-%20Report%202021.pdf>
- [5] K. R. Forti V., Balde C.P. E-waste statistics guidelines on classification reporting and indicators. [Online]. Available: <http://collections.unu.edu/view/UNU:6477>
- [6] Y. Xu, G. ke Yang, J. Luo, and J. He, "An electronic component recognition algorithm based on deep learning with a faster squeezeNet," *Mathematical Problems in Engineering*, vol. 2020, pp. 1–11, 2020.
- [7] S. Kumar, D. Yadav, H. Gupta, O. Verma, I. Ansari, and C. W. Ahn, "A novel yolov3 algorithm-based deep learning approach for waste segregation: Towards smart waste management," *Electronics*, vol. 10, p. 14, 12 2020.
- [8] DTI. Refind, & stena. march 23 new robot system extracts dangerous and valuable items from waste using artificial intelligence.
- [9] . Refind. Optical battery sorter 500. [Online]. Available: <https://www.refind.se/optical-battery-sorter-500>
- [10] H. Karbasi, A. Sanderson, A. Sharifi, and C. Pop, "Robotic sorting of used button cell batteries: Utilizing deep learning," in *2018 IEEE Conference on Technologies for Sustainability (SusTech)*, 2018, pp. 1–6.
- [11] S. Gundupalli, S. Hait, A. Thakur, and A. Trivedi, *Classification of Recyclables from E-Waste Stream Using Thermal Imaging-Based Technique*, 12 2018, pp. 67–78.
- [12] W. Sterkens, D. Diaz-Romero, T. Goedemé, W. Dewulf, and J. R. Peeters, "Detection and recognition of batteries on x-ray images of waste electrical and electronic equipment using deep learning," *Resources, Conservation and Recycling*, vol. 168, p. 105246, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0921344920305619>
- [13] S. Z. Valtchev and J. Wu, "Domain randomization for neural network classification," *Journal of Big Data*, vol. 8, no. 1, pp. 1–12, 2021.
- [14] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 23–30, 2017.
- [15] A. Loquercio, E. Kaufmann, R. Ranftl, A. Dosovitskiy, V. Koltun, and D. Scaramuzza, "Deep drone racing: From simulation to reality with domain randomization," *IEEE Transactions on Robotics*, vol. 36, no. 1, pp. 1–14, 2020.
- [16] J. Tremblay, A. Prakash, D. Acuna, M. Brophy, V. Jampani, C. Anil, T. To, E. Cameracci, S. Boochoon, and S. Birchfield, "Training deep networks with synthetic data: Bridging the reality gap by domain randomization," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1082–10828, 2018.
- [17] Y. Movshovitz-Attias, T. Kanade, and Y. Sheikh, "How useful is photo-realistic rendering for visual learning?" in *Computer Vision – ECCV 2016 Workshops*, G. Hua and H. Jégou, Eds. Cham: Springer International Publishing, 2016, pp. 202–217.
- [18] W. Chen, H. Wang, Y. Li, H. Su, Z. Wang, C. Tu, D. Lischinski, D. Cohen-Or, and B. Chen, "Synthesizing training images for boosting human 3d pose estimation," in *2016 Fourth International Conference on 3D Vision (3DV)*, 2016, pp. 479–488.
- [19] N. Zhang, J. Donahue, R. B. Girshick, and T. Darrell, "Part-based rcnns for fine-grained category detection," *ArXiv*, vol. abs/1407.3867, 2014.
- [20] D. Lin, X. Shen, C. Lu, and J. Jia, "Deep lac: Deep localization, alignment and classification for fine-grained recognition," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1666–1674.
- [21] X.-S. Wei, C.-W. Xie, J. Wu, and C. Shen, "Mask-cnn: Localizing parts and selecting descriptors for fine-grained bird species categorization," *Pattern Recognition*, vol. 76, pp. 704–714, 2018.
- [22] M. Jaderberg, K. Simonyan, A. Zisserman, et al., "Spatial transformer networks," *Advances in neural information processing systems*, vol. 28, pp. 2017–2025, 2015.
- [23] J. Fu, H. Zheng, and T. Mei, "Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4438–4446.
- [24] H. Zheng, J. Fu, T. Mei, and J. Luo, "Learning multi-attention convolutional neural network for fine-grained image recognition," 10 2017, pp. 5219–5227.
- [25] Z. Yang, T. Luo, D. Wang, Z. Hu, J. Gao, and L. Wang, "Learning to navigate for fine-grained classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 420–435.
- [26] X. He and Y. Peng, "Weakly supervised learning of part selection model with spatial constraints for fine-grained image classification," in *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [27] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view rgb-d object dataset," in *2011 IEEE International Conference on Robotics and Automation*, 2011, pp. 1817–1824.
- [28] M. Blum, J. T. Springenberg, J. Wülfing, and M. Riedmiller, "A learned feature descriptor for object recognition in rgb-d data," in *2012 IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1298–1303.
- [29] L. Zheng, Y. Yang, and Q. Tian, "Sift meets cnn: A decade survey of instance retrieval," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 5, pp. 1224–1244, 2017.
- [30] A. Eitel, J. T. Springenberg, L. Spinello, M. Riedmiller, and W. Burgard, "Multimodal deep learning for robust rgb-d object recognition," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 681–687.
- [31] B. O. Community, *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. [Online]. Available: <http://www.blender.org>
- [32] W. Schroeder, K. M. Martin, and W. E. Lorensen, *The visualization toolkit an object-oriented approach to 3D graphics*. Prentice-Hall, Inc., 1998.
- [33] A. Fedorov, R. Beichel, J. Kalpathy-Cramer, J. Finet, J.-C. Fillion-Robin, S. Pujol, C. Bauer, D. Jennings, F. Fennessy, M. Sonka, J. Buatti, S. Aylward, J. Miller, S. Pieper, and R. Kikinis, "3d slicer as an image computing platform for the quantitative imaging network," *Magnetic resonance imaging*, vol. 30, pp. 1323–41, 07 2012.
- [34] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 224–236.
- [35] B. Zhou, A. Khosla, L. A., A. Oliva, and A. Torralba, "Learning Deep Features for Discriminative Localization." *CVPR*, 2016.
- [36] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [37] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*. Springer, 2010, pp. 177–186.
- [38] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.
- [39] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," *CoRR*, vol. abs/1703.06907, 2017. [Online]. Available: <http://arxiv.org/abs/1703.06907>
- [40] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011.
- [41] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [43] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [44] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision. 2015," *arXiv preprint arXiv:1512.00567*, 2015.

- [45] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [46] S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, p. 64270–64277, 2018. [Online]. Available: <http://dx.doi.org/10.1109/ACCESS.2018.2877890>
- [47] M. Grandini, E. Bagli, and G. Visani, "Metrics for multi-class classification: an overview," 2020.
- [48] M. B. Shaikh and D. Chai, "Rgb-d data-based action recognition: A review," *Sensors*, vol. 21, no. 12, p. 4246, 2021.
- [49] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.