

# A Novel Incremental Learning Driven Instance Segmentation Framework to Recognize Highly Cluttered Instances of the Contraband Items

Taimur Hassan<sup>\*</sup>, *Member, IEEE*, Samet Akcay, Mohammed Bennamoun, *Senior Member, IEEE*, Salman Khan, and Naoufel Werghi, *Senior Member, IEEE*

**Abstract**—Screening cluttered and occluded contraband items from baggage X-ray scans is a cumbersome task even for the expert security staff. This paper presents a novel strategy that extends a conventional encoder-decoder architecture to perform instance-aware segmentation and extract merged instances of contraband items without using any additional sub-network or an object detector. The encoder-decoder network first performs conventional semantic segmentation and retrieves cluttered baggage items. The model then incrementally evolves during training to recognize individual instances using significantly reduced training batches. To avoid catastrophic forgetting, a novel objective function minimizes the network loss in each iteration by retaining the previously acquired knowledge while learning new class representations and resolving their complex structural interdependencies through Bayesian inference. A thorough evaluation of our framework on two publicly available X-ray datasets shows that it outperforms state-of-the-art methods, especially within the challenging cluttered scenarios, while achieving an optimal trade-off between detection accuracy and efficiency.

**Index Terms**—Baggage X-ray Scans, Semantic Segmentation, Instance Segmentation, Incremental Learning.

## I. INTRODUCTION

THE inspection of passenger's baggage, packages, and containers with X-ray scanners is nowadays a part of the standard checking measures in airports and any other public place where safety and security are of significant concern. This screening process is cumbersome, requiring the relentless attention of a human expert. Furthermore, it's vulnerable to human errors caused due to exhausting work schedules, lack of experience, and the concealed nature of the contraband items. Although object detection in color images has been a rigorously researched topic, its applicability to X-ray-based threat detection is somewhat limited. The primary reason

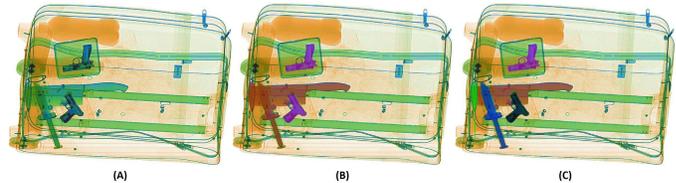


Figure 1: (A) An original X-ray scan from the SIXray dataset [1], (B) conventional semantic segmentation, and (C) instance-aware segmentation.

is the remarkably different X-ray imagery characteristics, where texture and appearance details are scarce compared to regular color images. An adequate system for such a critical application is expected to detect objects under high occlusions, in cluttered scenes, with large view-point variations and limited amounts of contraband data. Many researchers have developed supervised and unsupervised screening systems for detecting contraband items in X-ray images in response to these challenges. The most recent wave of these efforts employed deep learning models, particularly one-staged and two-staged object detectors such as RetinaNet [2], YOLO [3], and Faster R-CNN [4]. While these systems showed remarkable capacity for detecting isolated objects, their performance degrades in recognizing extremely cluttered, occluded, and overlapping items [5], [6]. Semantic segmentation models, due to their pixel-level recognition ability, can extract the extremely occluded contraband items from X-ray baggage scans [7]. With the integration of object context in the pixel classification, they have more potential to improve the threat detection accuracy [8]. By leveraging this capacity, some of the initial attempts employed the encoder-decoder-encoder topology for detecting suspicious items as anomalies [9]. However, semantic segmentation networks have an inherent limitation of detecting the individual instances of the overlapping items. For example, in Figure 1 (B), we can see that how a semantic segmentation network cannot recognize the overlapping *kitchen knife* and *chopper* individually. In such scenarios, these networks output only a single blob in which the information about individual item instances is lost. Detecting individual instances of the same threat category is, in fact, desirable in cases where we need to identify and locate each instance precisely (see the example in Figure 1-C, where the *kitchen knife* and the *chopper* instances have been extracted separately). Also, identifying individual items' instances is

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

This work is supported with a research fund from Khalifa University: Ref: CIRA-2019-047.

T. Hassan and N. Werghi are with the Center for Cyber-Physical Systems (C2PS), Department of Electrical Engineering and Computer Sciences, Khalifa University, Abu Dhabi, United Arab Emirates.

S. Akcay is with Intel R&D UK, United Kingdom.

M. Bennamoun is with the Department of Computer Science and Software Engineering, The University of Western Australia, Perth, Australia.

S. Khan is with the Mohamed Bin Zayed University of Artificial Intelligence, Abu Dhabi, United Arab Emirates.

<sup>\*</sup> Corresponding Author. Email: taimur.hassan@ku.ac.ae

vital in aviation baggage screening as some instances of the items are legal to carry within the baggage, whereas some instances are prohibited. For example, passengers can carry certain *drugs* and *bottles* in their luggage, but *addictive drugs* and *alcoholic drinks* are banned at airports [10]. Towards this end, Gaus et al. [6], [11] introduced an instance segmentation approach in their baggage threat detection system using Mask R-CNN [12]. However, the authors realized that conventional instance segmentation network requires extensive ground truth labeling and exhaustive training efforts, especially for the large-scale datasets, and there is a need to develop a framework that can effectively perform instance-aware segmentation to recognize the cluttered contraband items from the baggage X-ray imagery via incremental few-shot training.

## II. RELATED WORK

Existing solutions for contraband item detection based on X-ray imagery can be classified as traditional machine learning and deep learning methods. In this section, we shed light on the main approaches, and we refer the reader to the work of [13], and [14] for a detailed survey. In addition to this, this section also explores the recent advances in incremental learning to perform classification and segmentation tasks.

**A. Conventional Machine Learning Methods:** The initial methods developed for screening contraband items employ conventional machine learning. These solutions are either based on classification [15], detection [16] or the segmentation approaches [17]. Bastan et al. [18] used SURF features with Bag of Words (BoW) to identify suspicious objects. Instead of SURF, Kundegorski et al. [19] utilized FAST-SURF with BoW to classify prohibited baggage items. Other works involve Adaptive Sparse Representation [20], and Adapted Implicit Shape Model [21] to detect contraband data. Apart from this, Mery et al. [20] developed a framework that computes 3D feature points through the structure from motion and uses these features to classify contraband items from the X-ray imagery.

**B. Deep Learning Methods:** The most recent deep learning methods can be categorized either as supervised detection and segmentation approaches or as unsupervised adversarial learning schemes.

**1. Supervised Detection Strategies:** The majority of deep contraband item detection frameworks utilizes one-staged or two-staged object detectors such as YOLOv2 [22], RetinaNet [2] and Faster R-CNN [4]. Moreover, researchers have also utilized pre-trained models for the object classification within baggage X-ray scans [6], [23]. Zou et al. [24] utilized YOLOv2 [22] to detect *scissors*, *knives* and *bottles* from their local 1,104 synthetic X-ray images. Miao et al. [1] released the largest security inspection X-ray dataset (SIXray) that contains highly occluded and overlapping instances of contraband items such as *guns*, *knives*, *wrenches*, *pliers*, *scissors* and *hammers*. Furthermore, they presented a framework dubbed class-balanced hierarchical refinement (CHR) to recognize contraband items from the SIXray [1] dataset. More recently, Hassan et al. [25] presented Cascaded Structure Tensor (CST) framework that generates contours-driven bounding boxes of potentially prohibited items which are then classified using ResNet<sub>50</sub> [26].

**2. Supervised Segmentation Approaches:** Apart from solving the baggage threat recognition problem via deep object detection methods, many researchers utilized semantic and instance segmentation as a tool to effectively recognize suspicious baggage content [6], [8]. It is essential to note here that although we can fine-tune standard encoder-decoder networks for a large variety of semantic segmentation tasks, specific applications would be best be approached with customized models [27]. For example, to cope with object size variation and camera view changes in traffic and surveillance applications, Akilan et al. [28] proposed integrating residual feature fusions at early, middle and late stages in the encoder-decoder architecture (dubbed MvRF-CNN [28]). Similarly, driven by achieving the optimal trade-off between the segmentation accuracy and the computational model complexity, Wang et al. [29] coupled an encoder-decoder model and super-resolution construction scheme. Similarly, a multi-task attention network is proposed in [30] that coupled handcrafted features pipeline and an attention network to segment the object of interest [30]. Also, an adversarial domain adaptation scheme is proposed in [31] that employs a detection and segmentation (DS) model along with domain classifiers to learn target domain labels from the source domain synthetic data in a weakly supervised manner. In addition to this, Hassan et al. [7] proposed a contour instance segmentation strategy that segments the suspicious baggage content by analyzing the strength of the variation within their contours [7].

**3. Unsupervised Adversarial Learning:** Apart from supervised learning frameworks for detecting contraband items, Akcay et al. proposed GANomaly [9], and Skip-GANomaly [32] to derive the latent space representation of the contraband items in an adversarial manner to recognize them as anomalies within the baggage X-ray scans.

**C. Incremental Learning Strategies:** Incremental learning schemes have gained immense popularity in the context of deep learning for overcoming the need for excessive computational burden in re-training models with large-scale data, which might be difficult to obtain and prepare. However, developing an incremental learning scheme that overcomes catastrophic forgetting (the tendency of a deep learning model to drastically forget the prior knowledge while learning about new information) is also challenging. To address this, many researchers have proposed schemes involving knowledge distillation [33], gating [34], and indefinitely long term learning (iCaRL) [35]. Furthermore, Tian et al. [36] exploited the fact that knowledge representations exhibit complex relationships that cannot be learned through objective functions that assume independence of events. Cho et al. [37] advocated that good performing teachers do not necessarily produce good students due to the student network's limited capacity to cope with the teacher's growing knowledge. Lopez-Paz et al. [38] proposed the Gradient Episodic Memory (GEM) scheme, which uses episodic memories to hold a small set of examples from the prior learned tasks to avoid catastrophic forgetting. Apart from this, researchers have also proposed distillation-driven incremental learning strategies for performing the semantic segmentation tasks [39].

**D. Limitations of Existing Work:** The main limitations of the

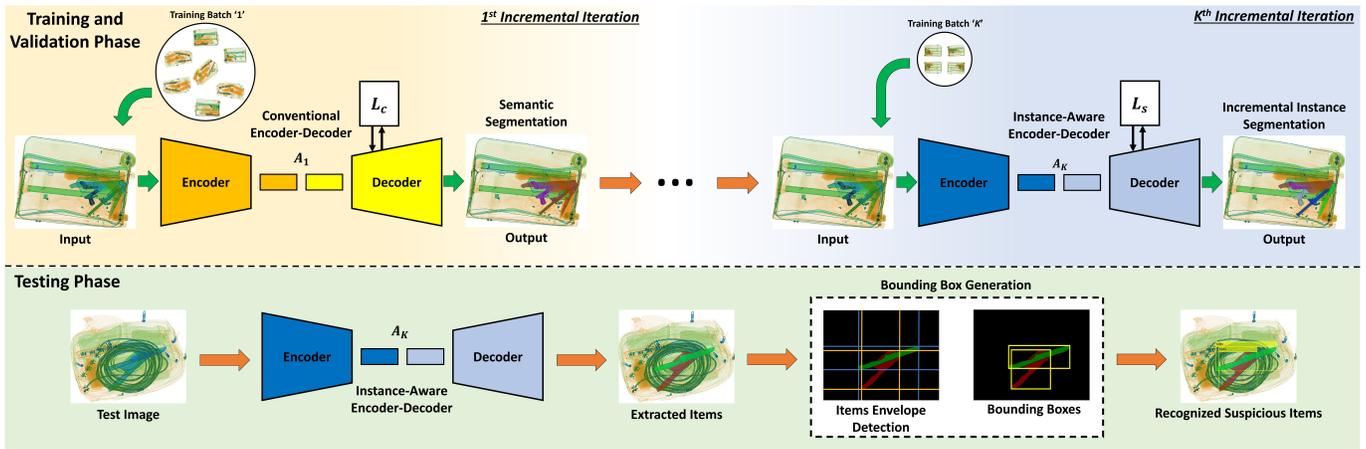


Figure 2: Block diagram of the proposed framework. We trained the proposed model incrementally to recognize cluttered instances of the contraband items. At each iteration,  $k = 1, \dots, K$ , the number of item instances that the system can recognize is incremented by one. At the inference stage, the model  $A_K$  (incrementally trained till  $K^{th}$  iteration to recognize up to  $K$  overlapped instances) is used for the instance-aware segmentation of the cluttered contraband items. More on the training details are in Section (IV-B).

existing approaches are their inadequate validation on single datasets or their application to simplistic scenarios within a very constrained environment. For instance, the problem of robustly detecting cluttered, occluded, and overlapping contraband items from the highly imbalanced datasets is still an open question to be addressed. The approaches proposed in [1], [25] and [7] handles such cases. However, they produce either low detection performance [1] or are subject to parameter tuning [25]. Apart from this, researchers have also utilized semantic segmentation networks to recognize suspicious baggage content via X-ray imagery [8]. Such models have improved the performance of the threat detection frameworks. However, they cannot distinguish between cluttered and overlapping instances of the same items (e.g., a *knife* overlaid on another *knife* as shown in Figure 1-B), which is often desirable in aviation screening, and for such cases, the semantic segmentation networks output a single blob of pixels representing only a single class label. To cater this, Gaus et al. [6] introduced the usage of Mask R-CNN [12] for baggage threat detection. However, the Mask R-CNN-based threat detection system presents limitations in extracting the cluttered contraband items because it relies on the region-based proposals that fail to detect cluttered objects correctly [6]. This limitation of Mask R-CNN [12] and other instance-aware segmentation networks will be further evidenced when employed in complex datasets such as SIXray [1], as described in Section V. Moreover, other approaches utilized encoder-decoder architectures and fully convolutional networks coupled with classification sub-networks or region of interest (ROI) voting to recognize multiple objects instances individually [8], [40]. However, these frameworks also produce a poor trade-off between detection accuracy and efficiency. On the other hand, instance segmentation frameworks require extensive bounding box and mask-level annotations [7], which are reasonably hectic, and resource-demanding to procure, especially for large-scale datasets, such as SIXray [1]. Also, training such networks requires an excessive amount of memory and computational resources. To alleviate these problems,

we propose an incremental learning-driven instance-aware segmentation approach, as discussed below.

**E. Contributions:** This paper proposes a novel scheme that utilizes incremental learning to make conventional semantic segmentation models instance-aware. The proposed method is simple and exhibits modest training efforts by requiring only a small batch of training samples to add more instances of a given suspicious item class. This strategy bypasses hectic annotation workflows as are necessary for training traditional instance segmentation frameworks while overcoming the excessive memory and computational requirements. The proposed framework also avoids catastrophic forgetting through an instance segmentation objective function that minimizes the network loss to retain knowledge about the previously learned classes while understanding new class representations and resolving their complex inter-dependencies. The unique characteristics of the proposed system are:

- A novel approach that extends conventional encoder-decoder networks to recognize individual instances of the contraband items from the X-ray scans.
- No requirement for an additional object detector, classification sub-network, or ROI voting to perform instance-aware segmentation.
- An incremental learning-driven instance segmentation framework that discriminates the overlapping and isolated suspicious item instances with only a few training examples.
- Robust to catastrophic forgetting due to its ability to resolve complex inter-dependencies between already learned and newly added suspicious items categories.

The rest of the paper is organized as follows: Section III discusses the proposed system. Section IV enlists the experimental plan. Section V presents the experimental results. Section VI contains a detailed discussion on the performance of the proposed system and Section VII presents concluding remarks.

### III. PROPOSED FRAMEWORK

Figure 2 depicts the block diagram of the proposed framework. This framework trains an encoder-decoder model to recognize up to  $K$  isolated and overlapped instances of a given class incrementally in  $K$  iterations. The first iteration reflects the ordinary semantic segmentation to extract different contraband items from the baggage X-ray scans. For this, we train the first instance of encoder-decoder dubbed  $A_1$  on a relatively large set of training images. Afterward, we make the encoder-decoder model instance-aware in each iteration by exposing it to the small training batches. For example, in the  $k^{th}$  iteration, we make the encoder-decoder  $A_k$  to recognize up to  $k$  instances of the same item by providing a different set of corresponding images. The final instance-aware segmentation model is obtained at the iteration  $K$ . In this process, the model is immunized to catastrophic forgetting by analyzing the complex relationships between previously learned and newly added suspicious item categories through the proposed loss function (see Eq. 2). Before exposing the details of our approach, we provide a brief overview of the incremental learning paradigm in the next section for completeness.

**A. Incremental Learning:** In a conventional incremental learning paradigm, the model is trained iteratively. At each iteration  $k$ , it performs  $C_W$ -class segmentation (or classification) task where  $C_W$  denotes the number of classes in the current iteration  $k$ . To learn this task, the model is given a set of  $\mathcal{D}$  training samples such that  $\mathcal{D} = \{\mathcal{D}_o, \mathcal{D}_n\}$ , where  $\mathcal{D}_o$  denotes the samples of old classes  $W_o$ , learned from iteration  $l$  to  $(k-1)$ , and  $\mathcal{D}_n$  denotes the samples of newly added classes ( $W_n$ ) to be learned in the current iteration  $k$ . The cumulative list of all the classes (both  $W_o$  and  $W_n$ ) is represented by  $W$ , i.e.,  $W = \{W_o, W_n\}$ . The network is also fed with the ground truth  $t = \{t^o, t^n\}$  of these training samples where  $t^o$  and  $t^n$  denote the ground truth for the samples of old classes and the new classes, respectively.  $t$  is normally represented in a one-hot encoding vector notation [41]. These training samples are passed as an input to the network for which it generates the output logits  $l$  in the last layer such that  $l = vf + \gamma$ , where  $f$  represents the feature vector,  $v$  represents the layer weights, and  $\gamma$  denotes the biasing factor. The logits  $l = \{l^o, l^n\}$  are the concatenation of the old logits  $l^o$  and the new logits  $l^n$ , generated from training the old classes and the newly added classes, respectively. These logits are then passed through the activation function (usually softmax) in the final layer of the CNN model to generate the final class probabilities, i.e.,  $p(l_{i,j}) = \frac{\exp(l_{i,j})}{\sum_{r=0}^{C_W-1} \exp(l_{i,r})}$ , where  $p(l_{i,j})$  denotes the probability of the  $i^{th}$  training sample being part of the  $j^{th}$  class.  $p(l_{i,j})$  in the above definition is known as a hard class probability of the logit  $l_{i,j}$ . Hard probabilities are generally recommended in traditional classification or segmentation task because they clearly discriminate the most expected class out of the rest. But in incremental learning, logits are scaled using the temperature constant ( $\tau$ ) to generate the soft target probabilities, i.e.,  $p(l_{i,j}^\tau) = \frac{\exp(l_{i,j}^\tau)}{\sum_{r=0}^{C_W-1} \exp(l_{i,r}^\tau)}$ , where  $l_{i,j}^\tau = l_{i,j}/\tau$ . Here,  $\tau$  is used to increase the degree of relaxation of the soft label by reducing the disparities between

classes probabilities. Practically, it is a hyper-parameter which is tuned for the sake of obtaining a better performing model [42].

**B. Semantic Segmentation:** The first iteration of the proposed framework relates to semantic segmentation, where we train the proposed contraband items extraction network (CIE-Net) to extract different contraband items from the baggage X-ray images. The prime objective of designing the proposed CIE-Net is to accurately extract the contraband items and their instances, even in overly cluttered scenarios. We utilize convolutional blocks (with ReLU activations and batch normalizations) to preserve coarser feature representations of the contraband items while simultaneously retaining their geometrical shapes through finer edge information. The blocks follow a hierarchical design to yield multi-scale representations of threat objects for superior mask-level extraction. Furthermore, we implant novel identity blocks within the encoder topology of the CIE-Net that further aids in preserving the object's geometrical characteristics regardless of the amount of clutter. The optimal values for the number of filters and kernel sizes are determined empirically after analyzing the similarly designed frameworks like PSPNet [43], and ResNet [26] to craft out the optimal design schematics for the CIE-Net.

The detailed architecture of CIE-Net is illustrated in Figure 3. Here, we can observe that the CIE-Net consists of an asymmetric encoder-decoder topology. The desired objects' contextual and geometrical features are preserved through the contextual preservation blocks (CPB), composed of cascaded convolution and batch normalization operations. CPB ensures that the network learns to discriminate the similar textured contraband items (even the cluttered ones) by tuning the network weights based upon categorical cross-entropy loss function ( $L_c$ ) in the first iteration, and the proposed instance segmentation loss function ( $L_s$ ) in the rest of the iterations. Moreover, to ensure that the network retains the finer shape representations of the contraband items, dedicated identity blocks (inspired by ResNet [26] scheme) have been added in the encoder part, where the finer representations (of the suspicious items) are fused with the decoder end via residual triggered skip-connections. Inspired by PSPNet [43], we also employ a custom hierarchical block (HB) to improve the performance of the CIE-Net further. HB uses variable pooling factors (determined empirically) to generate the multi-scale feature representations from the latent vector space to recognize the cluttered contraband items and their instances. The hierarchical decomposition and pooling factors are determined empirically to obtain the optimal contraband item extraction performance on grayscale and colored baggage X-ray scans. Like the proposed framework, the MvRF-CNN [28] also preserves the desired objects' geometrical information by fusing feature representations obtained across various network depths in a residual manner [28]. Similarly, to achieve better geometrical characteristics of the desired objects, the framework proposed in [29] couples a segmentation encoder-decoder model with the super-resolution construction scheme where the fine-grained structural features are derived through the affinity maps [29]. To have a precise idea of how the above model works to detect baggage threats from security

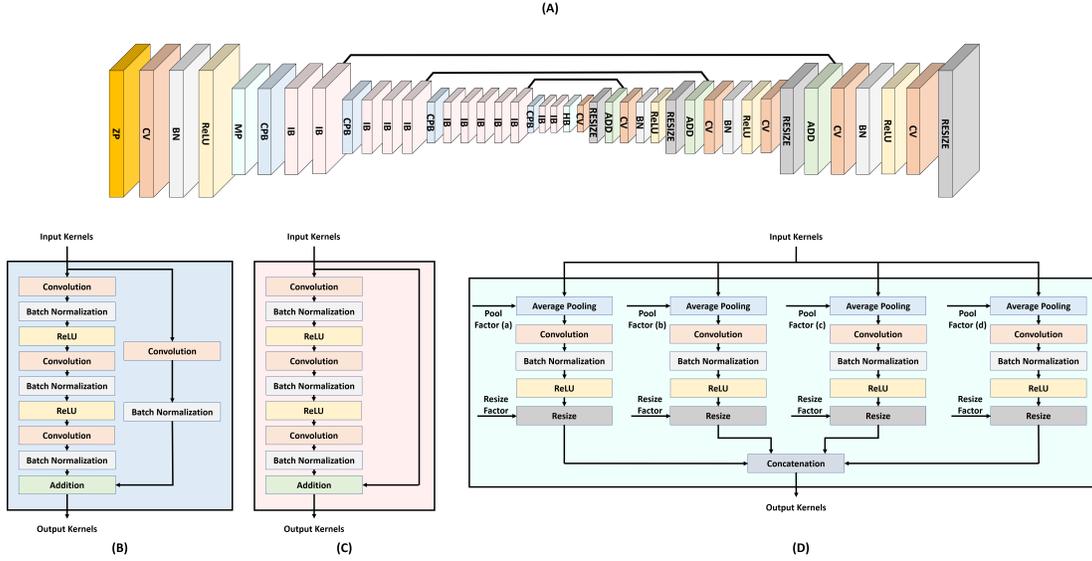


Figure 3: (A) CIE-Net architecture, (B) contextual preservation block (CPB), (C) identity block (IB), (D) hierarchical block (HB). Moreover, CV, BN, MP, and ZP in (A) denote the convolution, batch normalization, max pooling, and the zero-padding layer, respectively.

X-ray scans, we evaluated them both on the GDXray [44], SIXray [1], and the combined datasets. We also compared these scheme's performance with the proposed incremental instance segmentation framework (please see Table 4 for more details).

In the first incremental training iteration, CIE-Net optimizes the  $L_c$  function to discriminate between normal and suspicious items (in a semantic segmentation fashion):

$$L_c = -\frac{1}{N_t} \sum_{i=0}^{N_t-1} \sum_{j=0}^{C_W-1} t_{i,j} \log(p(l_{i,j})), \quad (1)$$

where  $C_W$  denotes the total number of classes for the current iteration,  $N_t$  represents the total number of samples in the training batch, for the current iteration,  $t_{i,j}$  is a binary value telling whether the  $i^{\text{th}}$  sample represents the  $j^{\text{th}}$  class or not, and  $p(l_{i,j})$  is the probability of the logit ( $l_{i,j}$ ) of  $i^{\text{th}}$  sample for the  $j^{\text{th}}$  class.

Here, we also want to highlight that the semantic segmentation network extracts isolated and merged suspicious items from the baggage X-ray scans in the first iteration. However, the network cannot differentiate between multiple instances of the same item (e.g., two or more *knives* or *guns* in a single scan, whether they are isolated or merged).

**C. Incremental Instance Segmentation:** We propose a novel instance segmentation framework that utilizes incremental learning to make conventional semantic segmentation networks instance-aware. Most of the instance-aware segmentation models employ object detectors, ROI voting, or separate classification sub-networks. However, such implications require additional overheads for preparing large-scale training data and excessive memory requirements. Contrary to this, our proposed scheme makes conventional encoder-decoder models instance-aware without needing any additional resources. Thanks to the incremental adaptation strategy, only a small-scale training batch is required in each iteration to learn about

multiple item instances in each scan, which drastically reduces the memory and computational requirements compared to the fine-tuning approaches. Furthermore, our framework has an in-built capacity to resist catastrophic forgetting through the proposed incorporation of the mutual information loss function, which analyzes the complex inter-dependencies between prior knowledge and newly learned information through Bayesian inference.

**1. Instance Segmentation Loss Function:** For instance-aware segmentation, we propose the following loss function.

$$L_s = \alpha_1 L_n + \alpha_2 L_o + \alpha_3 L_{mi}, \quad (2)$$

where  $\alpha_{\{1,2,3\}}$  denote the loss weights (determined empirically to be 0.2, 0.3, and 0.5).  $L_n$  minimizes the network loss for learning new instance categories, and  $L_o$  minimizes the distillation loss for retaining the prior learned knowledge (about segmenting the suspicious baggage items). Both  $L_n$  and  $L_o$  are widely used in continual learning frameworks to avoid catastrophic forgetting [42]. In the proposed framework,  $L_o$  is calculated through categorical cross-entropy loss, while  $L_n$  is calculated through KL divergence loss, as shown below:

$$L_o = -\frac{1}{N_{t_o}} \sum_{i=0}^{N_{t_o}-1} \sum_{j=0}^{C_{W_o}-1} t_{i,j}^o \log(p(l_{i,j}^o)), \quad (3)$$

$$L_n = \frac{1}{N_{t_n}} \sum_{i=0}^{N_{t_n}-1} \sum_{j=0}^{C_{W_n}-1} q(t_{i,j}^n) \log \left( \frac{q(t_{i,j}^n)}{p(l_{i,j}^n)} \right), \quad (4)$$

where  $N_{t_o}$  and  $C_{W_o}$  denote, respectively, the number of old training samples and the number of old classes (added in 1 to  $k-1$  iterations).  $N_{t_n}$  and  $C_{W_n}$  denote, respectively, the number of new training samples and the number of newly added categories (in the current  $k^{\text{th}}$  iteration).  $t_{i,j}^o$  and  $t_{i,j}^n$  represent, respectively, the ground truth for the training samples of the old and the new classes.  $p(l_{i,j}^o)$  is the predicted distribution of the scaled logits generated through the training

samples of the old classes.  $q(t_{i,j}^n)$  is the actual distribution generated from the true labels of the newly added classes,  $p(l_{i,j}^{n,\tau})$  represents the predicted distribution of the scaled logits generated through the training samples of the new classes.

$L_{mi}$  in Eq. (2) is the new proposed loss term, which we introduce to account for the inter-dependencies between old knowledge and newly learned information in our problem. More about the rationale and the description of this loss term is given in the next sub-section.

**2. Mutual Information Loss Function** The mutual information loss function ( $L_{mi}$ ) is based on the Bayesian inference that exploits the complex inter-dependencies between previously learned class representations (in iteration 1 to  $k-1$ ) through their respective training examples and the examples related to the newly stacked classes (in the current iteration  $k$ ).  $L_{mi}$  is expressed as follows:

$$L_{mi} = -\frac{1}{N_{t_o}} \sum_{i=0}^{N_{t_o}-1} \sum_{j=0}^{C_{W_o}-1} t_{i,j}^o \log(p(w_j | l_{i,j}^{o,\tau}, l_{i,j}^{n,\tau})), \quad (5)$$

where  $N_{t_o}$  denotes the total number of training examples,  $C_{W_o}$  denotes the total number of old classes ( $W_o$ ), and  $t_{i,j}^o$  the ground truth for the training samples representing the previously added classes (in iterations 1 to  $k-1$ ). The posterior probability  $p(w_j | l_{i,j}^{o,\tau}, l_{i,j}^{n,\tau})$  is defined as:

$$p(w_j | l_{i,j}^{o,\tau}, l_{i,j}^{n,\tau}) = \frac{p(l_{i,j}^{o,\tau}, l_{i,j}^{n,\tau} | w_j) \times p(w_j)}{\sum_{k=0}^{C_{W_o}-1} p(l_{i,k}^{o,\tau}, l_{i,k}^{n,\tau} | w_k) p(w_k)}. \quad (6)$$

It should be noted here that the evidence  $\sum_{k=0}^{C_{W_o}-1} p(l_{i,k}^{o,\tau}, l_{i,k}^{n,\tau} | w_k) p(w_k)$  in Eq. (6) is an optional term because it only normalizes the probability distribution  $p(w_j | l_{i,j}^{o,\tau}, l_{i,j}^{n,\tau})$ , so that the sum of probabilities for all the outcomes is 1.

The rationale of encompassing  $L_{mi}$  stems from the fact that older class representations (learned across the  $k-1$  iterations) and the newly learned categories (in the  $k^{th}$  iteration) are non-mutually exclusive. For example, a network trained to extract *knives* (particularly *kitchen knives*) in the first iteration should be aware of the contextual similarity between *kitchen knives* and *choppers* (which it learns in the second iteration) since both of them are different type of *knives*.

To the best of our knowledge, all the knowledge distillation and incremental learning solutions handle catastrophic forgetting by separately minimizing the network loss involved in learning the new tasks and maintaining the prior learned knowledge inferred from the previous model (or teacher) instance. But the frameworks, trained using these loss functions assume that both older and newly added class representations are independent of each other, leading towards compromised performance, especially in those scenarios when the incrementally learned information highly correlates with one another. In our approach, the additional loss function ( $L_{mi}$ ) integrates the relationship between prior learned and recently stacked classes through their training examples and exploits it via Bayesian inference to maximize the capacity of the incremental learning process of differentiating contraband item instances.

**D. Bounding Box Generation:** After extracting the suspicious items from the candidate scan, the bounding box for each

extracted item ( $\zeta$ ) is generated through a simple yet very effective scheme. We iterate over the mask of each extracted contraband item ( $\zeta$ ) within the candidate scan, where for each mask, we find its minimum and maximum row value. The minimum row value represents the minimum row index within the candidate scan, where the mask value is one. Similarly, the maximum row value represents the maximum row index (within the candidate scan), where the mask is 1. Afterward, we take the image transpose and repeat the same process to get the minimum and maximum column index required to generate (and fit) the bounding box. The mathematical expression of the whole scheme is as follows:

$$\langle y_{min}, y_{max} \rangle = \langle \underset{0 \leq u \leq M-1}{\operatorname{argmin}}(\zeta_u), \underset{0 \leq v \leq N-1}{\operatorname{argmax}}(\zeta_v) \rangle, \quad (7)$$

$$\langle x_{min}, x_{max} \rangle = \langle \underset{0 \leq u \leq M-1}{\operatorname{argmin}}(\zeta_u^T), \underset{0 \leq v \leq N-1}{\operatorname{argmax}}(\zeta_v^T) \rangle, \quad (8)$$

$$\beta_b = [x_{min}, y_{min}, x_{max} - x_{min}, y_{max} - y_{min}], \quad (9)$$

where  $u, v \in \mathbb{W}$ ,  $M$  and  $N$  denotes the width and height of  $\zeta$ , respectively, and  $\beta_b$  denotes the bounding box of the candidate contraband item (generated via its extracted mask).

#### IV. EXPERIMENTAL SETUP

This section reports the datasets, the training details, and the evaluation metrics (used in the evaluation and also in the comparative study).

**A. Datasets:** We evaluated the proposed framework on publicly available GDXray [44], SIXray [1], and the combined dataset (containing the scans from both GDXray [44] and SIXray [1] datasets). We report the detailed description of these datasets in the supplementary material (and in the source code repository<sup>1</sup>) due to space constraints.

**B. Incremental Training Details:** To incrementally train the proposed framework on the GDXray [44] dataset, we used a total of 788 scans (400 scans for extracting originally identified suspicious items and 388 scans for the locally identified items). However, for the SIXray [1] dataset, we used 80% of the scans for training and 20% for evaluation as per the dataset standard [1]. Note that the number of incremental training iterations depends on the number of cluttered item instances within each dataset. In the combined dataset, we have a total of 1,067,381 scans in which 27,750 scans (13,663 positives and 14,087 negatives) were used for training purposes, and the rest of 1,039,631 scans were used in the evaluations. Such a training split also ensures assessing the resistance of the proposed framework against class imbalance.

Moreover, in the first training iteration, we constrain the network with the  $L_c$  loss function to recognize different contraband items. Here, the proposed model performs conventional semantic segmentation to extract, for example, a *gun* and a *knife* contained within the candidate scan. However, it should be noted that the semantic segmentation model cannot recognize the overlapping instances of the same item, i.e., a *gun* overlaid on another *gun*. In such scenarios, the semantic

<sup>1</sup> The complete source code and its documentation is available at: <https://github.com/taimurhassan/inc-inst-seg>.

segmentation models will output a single blob of *gun*-labeled pixels.

To accurately recognize the individual overlapped instances of contraband items (e.g., two overlapping *guns*), we further train our model iteratively. In each incremental iteration, we stack new classes, representing individual instances of the contraband items. Through their respective training examples, we re-tune the proposed model to make it instance-aware. For example, in the second iteration, we train the proposed model to recognize at most two overlapped instances of any suspicious item (e.g., two instances of *guns*, two instances of *knives* etc.) by stacking two additional classes representing *gun* and *knife* instance. We, therefore, feed the network with a small batch of training examples (containing at most two overlapping instances), where the two overlapping suspicious items (e.g., two overlapping *guns*) are marked with two different class labels in the ground truth. The same process is repeated across all the iterations until we obtain  $K$ -instance aware segmentation model where  $K$  denotes the maximum overlapping instances of the same item within the dataset. In addition to passing training examples representing the newly stacked classes, we also pass a few examples representing the previous classes (added in the iterations 1 to  $k-1$ ). The set of samples used to train the proposed model at each iteration is significantly lesser than the amount of data that is required by its competitors [1], [5], [7], [25], i.e., it only uses around 20% of the total training data (defined as per the dataset standard), wherein each increment, about 10% examples are added to retain the knowledge of the previously learned categories.

The training is conducted on a machine with an Intel Core i7-9750H@2.6 GHz processor and 32 GB RAM with a single NVIDIA RTX 2080 Max-Q GPU, cuDNN v7.5, and a CUDA Toolkit v11.0.221. The CIE-Net is implemented using TensorFlow 2.1.0 with Keras 2.3.0 on the Anaconda platform using Python 3.7.9. In the first iteration, the training consisted of 20 epochs, whereas the subsequent iterations took ten epochs with ADADELTA [45] optimizer. Moreover, the exact number of learnable and non-learnable parameters in CIE-Net varies in each iteration. Still, on average, they are roughly around 31.4M and 61.3K, respectively. The detailed model architecture is available in the codebase repository<sup>1</sup>.

We also tested the proposed framework’s applicability on the RGB data by evaluating it on the Microsoft COCO dataset [46]. Since the experiments on COCO dataset [46] do not relate to our proposed study, we report them in the supplementary material of this paper.

**C. Evaluation Metrics:** The proposed framework has been evaluated using the pixel-level recall, precision, intersection-over-union (IoU), dice coefficient (DC), ROC curves, box-level and mask-level mean average precision ( $\mu_{ap}$ ) computed using  $\text{IoU} \geq 0.5$  ( $\mu_{ap}^{b:50}$  and  $\mu_{ap}^{m:50}$ ),  $\text{IoU} \geq 0.75$  ( $\mu_{ap}^{b:75}$  and  $\mu_{ap}^{m:75}$ ), and  $\text{IoU} = 0.5 : 0.05 : 0.95$  ( $\mu_{ap}^b$  and  $\mu_{ap}^m$ ), respectively.

## V. RESULTS

This section reports a thorough evaluation of the proposed framework for extracting and recognizing the contraband items. The purpose of these experiments is two-fold: 1) comparing the performance of our instance segmentation model

Table 1: Evaluation of the different segmentation models on the SIXray (S) [1], GDXray (G) [44] and Combined (C) dataset. Bold indicates the best performance.

Model	IoU			DC		
	S	G	C	S	G	C
CIE-Net	<b>0.6883</b>	0.7723	<b>0.5861</b>	<b>0.8153</b>	0.8715	<b>0.7390</b>
CIE-R-Net	0.6702	<b>0.7852</b>	0.5749	0.8025	<b>0.8796</b>	0.7300
PSPNet	0.6641	0.7694	0.5728	0.7981	0.8696	0.7283
SegNet	0.6559	0.7463	0.5640	0.7921	0.8547	0.7212
U-Net	0.6434	0.7384	0.5514	0.7830	0.8495	0.7108
FCN-8	0.5792	0.6431	0.4527	0.6973	0.7827	0.6232
FCN-32	0.5084	0.6246	0.3931	0.6740	0.7689	0.5643

(CIE-Net) with other state-of-the-art models [12], [47]–[49], and 2) comparing the overall performance of our framework for baggage threat detection with other competitive systems [1], [7], [25], [50]. At first, we conducted an ablative analysis to assess the performance of different state-of-the-art encoder-decoder and fully convolutional models in our framework. We also conduct empirical experimentation to study the effect of the temperature constant ( $\tau$ ) and the effect of utilizing different knowledge distillation loss functions for incremental instance segmentation. Then, we present the detailed evaluation results of the proposed framework on both GDXray and SIXray datasets in Section V-B and Section V-C, respectively. Afterward, we report, in Section V-D, the experimentation conducted on the combined datasets.

**A. Ablation Study:** We conducted an ablation study to investigate: 1) The optimal choice of the segmentation model; 2) The effect of the temperature constant ( $\tau$ ); 3) The effects of employing different knowledge distillation loss functions in the incremental instance segmentation. Apart from this, we also conducted rigorous ablation experiments to evaluate the parametric effects of the CIE-Net and its custom CPB, IB, and HB blocks. Due to space constraints, these parametric evaluations are reported within the supplementary material of the article.

**1. Choice of Segmentation Model:** In this study, we compared the performance of several state-of-the-art semantic segmentation models, including PSPNet [43], SegNet [51], U-Net [52], FCN-8 and FCN-32 [53] with our proposed CIE-Net model for the extraction of isolated and overlapping contraband items and their instances depicted within the grayscale and colored baggage X-ray scans. We further want to notify that to fairly compare all the models, we have trained them incrementally using the proposed  $L_s$  loss function where each model, including the CIE-Net model, was implemented using ResNet<sub>101</sub> [26]. We dubbed this CIE-Net variant as CIE-R-Net to differentiate it from the CIE-Net build with our custom backbone.

The comparison results are reported in Table 1, where we can see that the proposed CIE-Net produced the best performance in terms of both IoU and DC metrics for the SIXray [1], GDXray [44], and the combined datasets.

Moreover, Figure 4 depicts a qualitative comparison showing segmentation results on samples from the SIXray and GDXray dataset. We can observe here that the CIE-Net produces better extraction results, especially for the examples in Figure 4 (A),

Table 2: Effects of varying the temperature parameter  $\tau$  (in terms of IoU).

$\tau$	GDXray	SIXray	Combined
0.1	0.4462	0.4106	0.2614
0.2	0.5013	0.4731	0.3053
0.5	0.6425	0.5632	0.3987
1	0.7341	0.6482	0.5014
1.5	0.7524	<b>0.6883</b>	0.5659
2	<b>0.7723</b>	0.6697	<b>0.5861</b>
2.5	0.7214	0.6021	0.5543
3	0.6642	0.5364	0.4471

(AJ), (AQ) and (AX). This better performance emanates from integrating the CPB, IB, and HB blocks in our model as showcased through rigorous parametric evaluations discussed in the supplementary material. Also, such synergy allows better extraction of contraband items by retaining global contextual information about the contraband items, even at the sparsest level of decomposition, while integrating finer features from the consecutive encoder part through the skip-connections.

**2. Effects of the Temperature Parameter:** In this experiment, we varied  $\tau$  from 0.1 to 3 and measured its effects on the segmentation performance for GDXray, SIXray, and the combined datasets. The results, depicted in Table 2, indicate  $\tau = 2$  and  $\tau = 1.5$  as the best values for the GDXray and the SIXray datasets, respectively.  $\tau = 2$  also yields the highest performance on the combined dataset. These results suggested framing the optimal values of  $\tau$  within the range [1.5, 2].

**3. Knowledge Distillation Loss Function:** This objective of this ablation study is to compare  $L_{mi}$  function with other state-of-the-art knowledge distillation loss functions, such as Output Distillation Loss ( $L_{od}$ ) [54], Modified Deep Model Consolidation [55] Loss ( $L_{ds}$ ) (proposed in [39]), Similarity-Preserving Knowledge Distillation Loss ( $L_{sp}$ ) [56], and Joint Classification and Distillation Loss ( $L_{cd}$ ) [35], in our framework. The comparison was made by switching the  $L_{mi}$  term in Eq. 2 with these distillation loss functions.

In what comes next, we denote by  $A_{k-1}$  and  $A_k$ , the models trained in the previous iteration (from 1 to  $k-1$ ), and in the current iteration  $k$ , respectively,  $N_{t_o}$  denotes the total number of training examples belonging to the previously learned classes,  $X_i^o$ ,  $i = 1 : N_{t_o}$ , denotes an old training sample,  $C_{W_o}$  denotes the total number of old classes, and  $F$  represents the Frobenius norm. Moreover,  $L_{od}$  minimizes the cross-entropy loss between the prediction of  $A_{k-1}$  and  $A_k$  and is expressed below:

$$L_{od} = \frac{1}{N_{t_o}} \sum_{i=0}^{N_{t_o}-1} \sum_{j=0}^{C_{W_o}-1} (p(l_{i,j}^{o,\tau})_{A_{k-1}}) \log(p(l_{i,j}^{o,\tau})_{A_k}), \quad (10)$$

$L_{ds}$  minimizes the disparities between the latent space feature representation of  $A_{k-1}$  and  $A_k$  and defined as:

$$L_{ds} = \frac{1}{N_{t_o}} \sum_{i=0}^{N_{t_o}-1} \|\mathcal{E}_{k-1}(X_i^o) - \mathcal{E}_k(X_i^o)\|_F^2, \quad (11)$$

where  $\mathcal{E}_{k-1}$  and  $\mathcal{E}_k$  are the latent space vectors related to  $A_{k-1}$  and  $A_k$ , respectively.

Table 3: Comparison of  $L_{mi}$  with state-of-the-art knowledge distillation loss functions in terms of IoU. To ensure fairness, we used CIE-Net with all the loss functions.

Loss Functions	GDXray [44]	SIXray [1]	Combined
$L_{mi}$	<b>0.7723</b>	<b>0.6883</b>	<b>0.5861</b>
$L_{sp}$ [56]	0.7504	0.6734	0.5480
$L_{od}$ [54]	0.7349	0.6162	0.5018
$L_{ds}$ [39]	0.7421	0.6395	0.5237
$L_{cd}$ [35]	0.6052	0.4793	0.2746

$L_{sp}$  minimizes the disparities between the activation similarity matrices ( $\mathcal{S}$ ) [56], and expressed as:

$$L_{sp} = \frac{1}{N_{t_o}} \sum_{i=0}^{N_{t_o}-1} \|\mathcal{S}_{k-1}(X_i^o) - \mathcal{S}_k(X_i^o)\|_F^2. \quad (12)$$

The joint classification and distillation loss  $L_{cd}$ , proposed in iCaRL [35], is expressed as follows:

$$L_{cd} = \mathcal{L}_{CE}(t_{i,j}^n, l_{i,j}^{n,\tau}) + \mathcal{L}_{CE}(t_{i,j}^o, l_{i,j}^{o,\tau}), \quad (13)$$

where  $\mathcal{L}_{CE}$  is the standard cross-entropy loss function and the other terms are as previously defined in Eq. (3) and (4). Note that unlike the previous knowledge distillation loss functions, which are plugged as a replacement to  $L_{mi}$ ,  $L_{cd}$  is used as a replacement of  $L_s$  in Eq. 2. This is because  $L_{cd}$  minimizes both the loss for learning new class representations and the distillation loss for retaining the previously learned classes.

The comparison of the loss functions is reported in Table 3 in term of IoU score where we can see that the proposed framework achieves 2.83%, 2.16%, and 6.50% performance improvements over the second-best  $L_{sp}$  [56] on GDXray [44], SIXray [1], and the combined dataset, respectively. These improvements emanate because of the synergy between  $L_n$ ,  $L_o$ , and  $L_{mi}$  that not only retains the prior knowledge while learning new classes but also enables the network to analyze the mutual relationships between the knowledge representations of the old and the new instances via Bayesian inference, unlike its competitors, that mostly rely on the spatial [54] and contextual [56] differences between knowledge representations.

**B. Evaluations on GDXray Dataset:** The CIE-Net was trained for two iterations on GDXray as this dataset contains at most two overlapping instances of the same contraband item. Table 5 shows the performance comparison against the state-of-the-art schemes. We can observe that our framework achieves 4.08% and 28.39% better performance than the second-best HTC [48] and the YOLACT [49], respectively, in terms of  $\mu_{ap}^m$ . Furthermore, it outperforms the second-best performing HTC [48] by 2.13% in terms of  $\mu_{ap}^b$ . However, for  $\mu_{ap}^{b:50}$ , the best performance is achieved by the original TST [7] (dubbed TST<sub>o</sub>) from which the proposed framework lags by 11.53%. However, this is an unfair comparison since TST [7] is trained conventionally using the large-scale well-annotated training data. In contrast, the proposed framework is trained incrementally on small-scale training batches. Moreover, under fair comparison with the incremental TST [7] scheme, dubbed TST- $L_s$ , the proposed CIE-Net is leading by 3.63%. Apart from this, the CIE-Net performance is further evaluated through the ROC curves, as shown in Figure 6 (a). These

Table 4: Comparison of the proposed framework with state-of-the-art solutions for extracting baggage threats. Bold indicates the best performance, while the second-best scores are underlined.

Metric	Method	GDxRay	SIXray	Combined
IoU	Proposed	<b>0.7723</b>	<b>0.6883</b>	<b>0.5861</b>
	MS RCNN [47]	0.7201	0.6484	0.5482
	Mask RCNN [12]	0.7098	0.6381	0.5243
	HTC [48]	0.7364	<u>0.6559</u>	<u>0.5804</u>
	YOACT [49]	0.7089	0.6110	0.4937
	DSRL [29]	<u>0.7421</u>	0.6542	0.5709
	MvRF-CNN [28]	<u>0.6982</u>	0.6016	0.4918
	TST- $L_s$ [7]	0.6851	0.5874	0.4285
DC	Proposed	<b>0.8715</b>	<b>0.8153</b>	<b>0.7390</b>
	MS RCNN [47]	0.8372	0.7867	0.7081
	Mask RCNN [12]	0.8302	0.7790	0.6879
	HTC [48]	0.8481	<u>0.7921</u>	<u>0.7344</u>
	YOACT [49]	0.8296	0.7585	0.6610
	DSRL [29]	<u>0.8519</u>	0.7909	0.7268
	MvRF-CNN [28]	<u>0.8222</u>	0.7512	0.6593
	TST- $L_s$ [7]	0.8131	0.7400	0.5999
Recall	Proposed	<b>0.8643</b>	<b>0.8057</b>	<b>0.7391</b>
	MS RCNN [47]	0.8238	0.7613	0.6846
	Mask RCNN [12]	0.8183	0.7542	0.6653
	HTC [48]	0.8392	<u>0.7736</u>	<u>0.7284</u>
	YOACT [49]	0.8195	0.7461	0.6548
	DSRL [29]	<u>0.8407</u>	0.7705	0.7173
	MvRF-CNN [28]	0.8196	0.7344	0.6419
	TST- $L_s$ [7]	0.8092	0.7269	0.5764
Precision	Proposed	<b>0.8952</b>	<b>0.8348</b>	<b>0.7401</b>
	MS RCNN [47]	0.8564	0.8153	0.7269
	Mask RCNN [12]	0.8439	0.8072	0.7154
	HTC [48]	0.8607	<u>0.8236</u>	<u>0.7318</u>
	YOACT [49]	0.8353	0.7669	0.6703
	DSRL [29]	<u>0.8736</u>	0.8125	0.7311
	MvRF-CNN [28]	0.8245	0.7801	0.6786
	TST- $L_s$ [7]	0.8173	0.7614	0.6256

curves are generated considering the pixel-level recognition, i.e., the pixel for each item (along with their instances) are treated as one and the rest of the pixels as zero (a typical binary classification). We can observe that the instance-aware CIE-Net achieved the minimum AUC score of 0.9818 for extracting *razors*. Due to space constraints, we report the detailed AUC score for each item (for all the datasets) within the source code repository<sup>1</sup>.

Moreover, we also compared the performance of the proposed CIE-Net against the state-of-the-art semantic and instance segmentation frameworks. The results are reported in Table 4, where we can see that on GDxRay, in terms of IoU, CIE-Net achieves 3.91% improvements over the DSRL [29] framework. Similarly, it outperforms HTC [48] by 4.64%.

In addition to this, we fairly compared the proposed framework with TST [7] by incrementally training it using the same experimental protocols and the proposed  $L_s$  function, where the proposed framework achieves 11.29% superior results, in terms of IoU, as evident from Table 4. The degradation in the TST's performance stems from the fact that during incremental training, it is more susceptible to forgetting the prior learned categories while adapting to new class representations since it employs a contour-driven strategy towards recognizing contraband items [7].

Moreover, the performance of CIE-Net on the GDxRay dataset is further analyzed through visual examples, as shown in

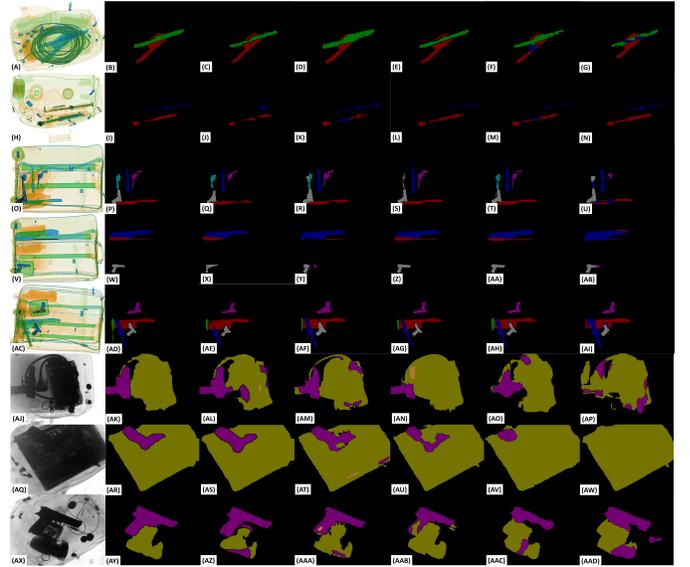


Figure 4: Extraction of contraband items (and their instances) using different segmentation models. From left: Original X-ray scan, CIE-Net, PSPNet [43], SegNet [51], U-Net [52], FCN-8, and FCN-32 [53]. Zoom-in for better visualization.

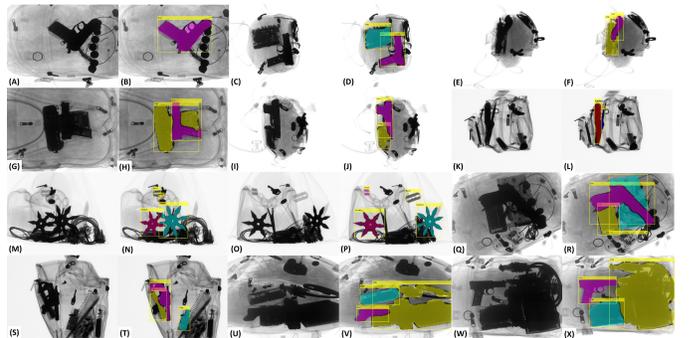


Figure 5: GDxRay [44]: Examples of occluded and overlapping items detection. Please zoom-in for better visualization.

Figure 5. The GDxRay contains at most two overlapping instances of the same items, e.g., see Figure 5 (N, L, P, R, V, and X). Here, we can appreciate the extraction performance of CIE-Net by observing two extracted occluded *knives* in (L) and occluded *shuriken* in (N, P). We can also observe how accurately the low-intensity *razors* have been segmented by the CIE-Net in Figure 5 (N, P).

**C. Evaluations on SIXray Dataset:** For the SIXray dataset, the training was conducted for six iterations since there are at most six instances of the same item in this dataset. Table 5 shows the model's comparison against the state-of-the-art instance segmentation algorithms. CIE-Net achieves 5.63% improvements in terms of  $\mu_{ap}^m$  against the second-best HTC [48] and 30.03% higher than the least good performing YOACT [49]. It also achieves 5.31% superior results than the existing solutions in terms of  $\mu_{ap}^b$ . For  $\mu_{ap}^{b:50}$ , the CIE-Net comes third after the original CST [25] (dubbed  $CST_o$ ) and the original TST [7] (dubbed  $TST_o$ ) scheme. However, this comparison is unfair, and the increased performance of  $CST_o$  [25] and  $TST_o$  [7] here emanates from the conventional

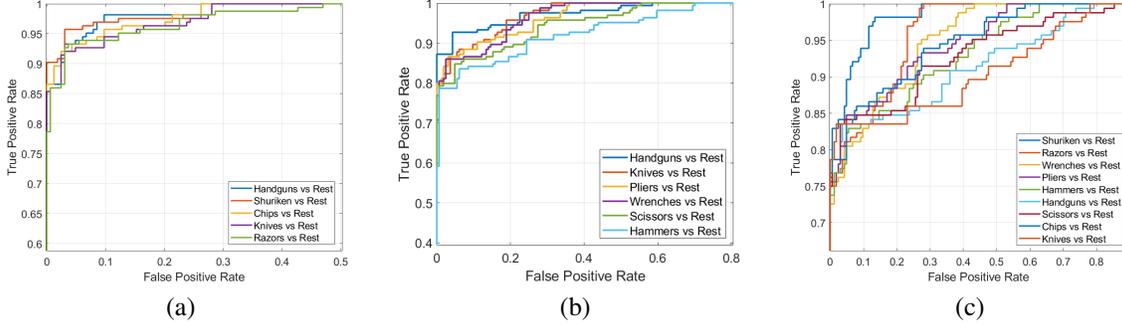


Figure 6: Performance evaluation of CIE-Net in terms of ROC for extracting contraband items from (a) GDXray dataset, (b) SIXray dataset, and (c) the combined dataset.

Table 5: Comparison of the proposed framework with state-of-the-art solutions for extracting contraband items. Bold indicates the best scores, while ‘-’ means that the metric is not computed.

D	M	$\mu_{ap}^m$	$\mu_{ap}^{m:50}$	$\mu_{ap}^{m:75}$	$\mu_{ap}^b$	$\mu_{ap}^{b:50}$	$\mu_{ap}^{m:75}$	
G	PF	<b>0.5068</b>	<b>0.7902</b>	<b>0.5006</b>	<b>0.6101</b>	0.8556	<b>0.6462</b>	
	MSR	0.4584	0.7283	0.4986	0.5564	0.8091	0.6033	
	MR	0.4311	0.7194	0.4893	0.5282	0.7833	0.5842	
	HTC	0.4861	0.7706	0.4997	0.5971	0.8314	0.6324	
	YT	0.3629	0.6518	0.3794	0.4852	0.7478	0.5491	
	CST <sub>o</sub> *	-	-	-	-	0.9343	-	
	TST <sub>o</sub> *	-	-	-	-	<b>0.9672</b>	-	
	TST <sub>i</sub>	-	-	-	-	0.8245	-	
	CST <sub>i</sub>	-	-	-	-	0.8169	-	
	TSD*	-	-	-	-	0.9162	-	
	S	PF	<b>0.4795</b>	<b>0.6893</b>	<b>0.4872</b>	<b>0.5367</b>	0.7653	<b>0.5374</b>
		MSR	0.4017	0.6347	0.4063	0.4653	0.6756	0.4782
MR		0.3654	0.5973	0.3592	0.4182	0.6326	0.4067	
HTC		0.4525	0.6629	0.4538	0.5082	0.7384	0.5021	
YT		0.3355	0.5632	0.3190	0.3811	0.6237	0.3643	
CST <sub>o</sub> *		-	-	-	-	<b>0.9595</b>	-	
TST <sub>o</sub> *		-	-	-	-	0.9516	-	
TST <sub>i</sub>		-	-	-	-	0.7248	-	
CST <sub>i</sub>		-	-	-	-	0.7351	-	
TSD*		-	-	-	-	0.6457	-	
CHR		-	-	-	-	0.5760	-	
C		PF	<b>0.4059</b>	<b>0.6249</b>	<b>0.4153</b>	<b>0.4862</b>	<b>0.7249</b>	<b>0.4983</b>
	MSR	0.3591	0.5986	0.3865	0.4023	0.6298	0.4572	
	MR	0.3129	0.5542	0.3301	0.3627	0.5983	0.3821	
	HTC	0.4023	0.6173	0.4102	0.4752	0.7203	0.4859	
	YT	0.3098	0.5286	0.3123	0.3561	0.5937	0.3696	
	TST <sub>i</sub>	-	-	-	-	0.6718	-	
	CST <sub>i</sub>	-	-	-	-	0.6526	-	

Abbreviations: D: Dataset, G: GDXray [44], S: SIXray [1], C: Combined Dataset, M: Methods, PF: Proposed Framework, MSR: Mask Scoring R-CNN [47], MR: Mask R-CNN [12], and YT: YOLACT [49]. Moreover, ‘\*’ indicates unfair comparison.

fine-tuning strategy, which utilizes the whole training dataset. Under fair comparison with incremental TST [7] (dubbed TST<sub>i</sub>) and CST (dubbed CST<sub>i</sub>), the CIE-Net is leading by 5.29% and 3.94%, respectively. Apart from this, the CIE-Net performance on SIXray is further evaluated through the ROC curves shown in Figure 6 (b). Here, we can observe that the proposed framework achieves the best AUC score for extracting the *handguns*. In addition to this, the segmentation performance of our framework can be analyzed through the mean IoU score in Table 4, showing the best score of 0.6883, leading the second-best HTC [48] by 4.70%.

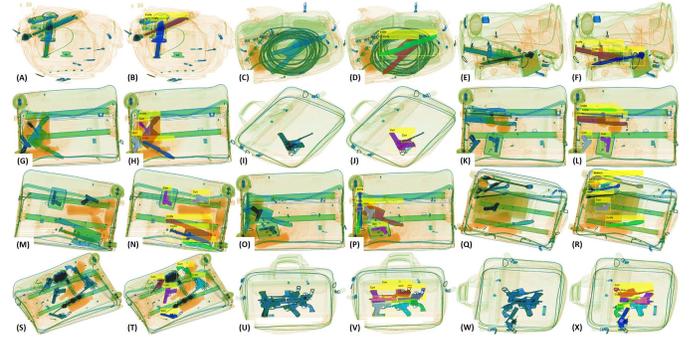


Figure 7: SIXray [1]: Examples of occluded and overlapping objects detection. Please zoom-in for better visualization..

In Figure 7, we report the qualitative evaluation showcasing examples of successfully extracted overlapping items, e.g., two items (B, D, F, H) and three items (N, P, R) and up to six items (V, X). In these examples, we can appreciate the potential of the instance-aware CIE-Net in accurately recognizing the extremely merged items, e.g., an instance of *guns* in Figure 7 (J, V, and X).

**D. Evaluations on Combined Dataset:** We also evaluated the proposed framework on the combined dataset. The results on the combined dataset are reported in Table 4 and 5. From Table 5, we can observe that CIE-Net achieved the best  $\mu_{ap}^{b:50}$  performance of 0.7249, outperforming the second-best framework by 0.6345%. Furthermore, we can also notice the performance gain of 23.67% over YOLACT [49] in terms of  $\mu_{ap}^m$ . Moreover, in terms of recall and precision, the CIE-Net is outperforming the second-best framework by 1.44%, and 1.12%, respectively (see Table 4).

In addition to this, Figure 6 (c) further depicts the ROC performance of instance-aware CIE-Net for extracting contraband items. Here, we can see that the minimum score is obtained for *knives* and *handguns* (i.e., AUC of 0.9133 and 0.9212, respectively).

Figure 8 showcases some qualitative examples derived from the combined dataset, which illustrates the capacity of CIE-Net for extracting instances of overlapped items despite the large differences of the scan properties in GDXray and SIXray datasets. In Figure 8 (F), we can observe how effectively the *razor* is extracted in such a cluttered scenario. Figure

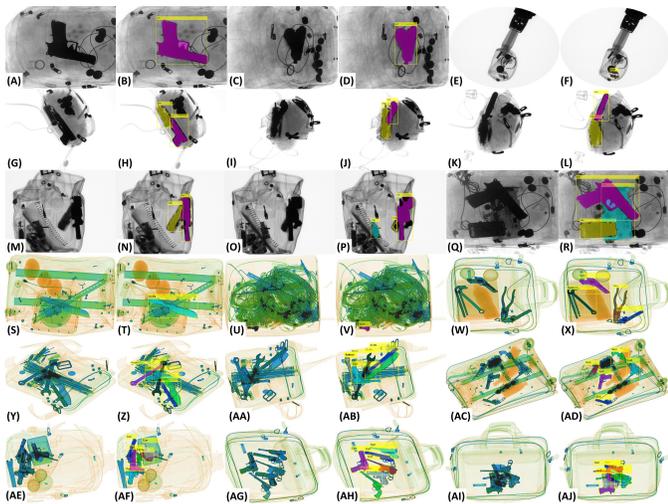


Figure 8: Examples of occluded and overlapping objects detection on combined dataset. Zoom-in for better visualization.

8 (N, R) depicts examples whereby our framework robustly differentiated between merged *gun* and *chip* instances. Figure 8 (T) depicts a reasonable extraction of the occluded *knife*. The performance of the CIE-Net can also be appreciated on more highly challenging scans such as (V), where a *gun* has been extracted from an extremely cluttered environment, (AB) in which two overlapping *wrenches*, two overlapping *knives* and a barely visible *gun* have been recognized, (AF) and (AJ) from which six extremely overlapping *guns* are accurately extracted. In Figure 8 (AF, AJ), in particular, we can appreciate the capacity of CIE-Net in accurately recognizing six instances of *guns* under extreme occlusion.

**E. Comparison of Run-time Performance:** Apart from evaluating the proposed scheme’s detection performance, we also analyzed its run-time performance and compared it with state-of-the-art methods. The comparison is reported in Table 6. Here, we can see that the proposed CIE-Net lags behind the state-of-the-art frameworks in terms of efficiency. This is due to the design choice of CIE-Net to focus more on accurately extracting the contraband items rather than achieving efficiency.

Due to this, the CIE-Net is slower than the other lightweight models like YOLOv3 [57], and CST [25]. However, we also want to highlight that the proposed framework is an instance segmentation scheme (unlike region-based YOLOv3 [57] and contour-based CST [25] detectors), and it gives the best trade-off between contraband items extraction (see Table 5) and run-time performance (see Table 6).

**F. Failure Cases:** Although the proposed framework achieves remarkable performance towards extracting overlapping contraband items (and their instances), as evident from Table 4, and 5, there are some cases where the CIE-Net turns out to be limited, especially on the negative SIXray scans (see pairs (A, B), (C, D), (E, F), (K, L) and (M, N) in Figure 9), producing pixel-level false positives and false negatives due to spatial and contextual similarity between the normal and suspicious baggage content within the X-ray scans. False positives are produced when the background regions (within the candidate

Table 6: Comparison of the run-time performance. The scores here represent the mean inference time of the two datasets. Bold indicates the best performance while the second-best performance is underlined.

Method	Time Performance (sec)
YOLOv3 [57]	<b>0.023</b>
CST [25]	<b>0.023</b>
RetinaNet [2]	<u>0.033</u>
YOLCAT [49]	0.036
CIE-Net (Proposed)	0.072
Mask R-CNN [12]	0.141
MS R-CNN [47]	0.156
HTC [48]	0.311

scan) are misclassified as threatening items by the proposed framework as shown in Figure 9-B, D, F, L, and N. Moreover, false negatives are generated when the region of the contraband item is misclassified as background. For example, see the missed portion of *shuriken* in Figure 9 (X). Apart from this, in some cluttered cases, the proposed CIE-Net produced over-segmentation results by confusing between different instances of the suspicious items (as shown in Figure 9-H, P, R, T, V, and Z). Although all these types of failures were seen rarely during the experimentation, they can be remedied through postprocessing schemes such as blob filtering, region-opening, and region-filling schemes.

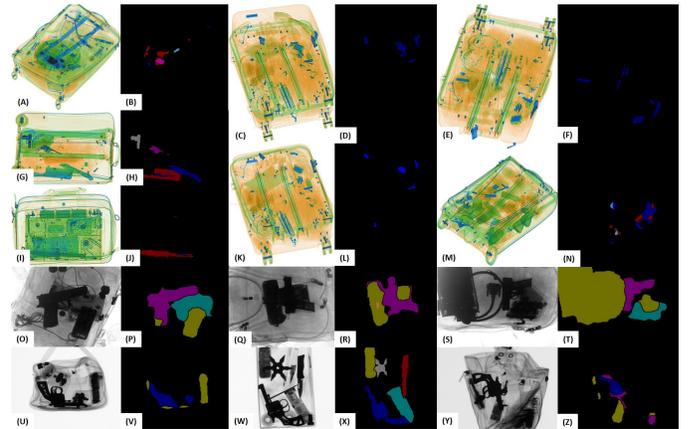


Figure 9: Failure cases in GDXray and SIXray datasets. Blue in (B, D, F, H, J, L, N) and red color in (B, H, J, N, R, X) represent *knives*. White in (H), magenta in (B, H, P, R, T, Z), cyan in (B, N), and blue color in (V, X, Z) represent *handguns*. Yellow and cyan color in (P, R, T, V, X, Z) depicts *chips*. White color in (V) represents the *shuriken*. Zoom-in for better visualization.

## VI. DISCUSSION

An overview of the results in Tables 4 and 5 convey that the proposed CIE-Net, employed within the incremental instance segmentation framework, shows neat performance improvement over standard models such as Mask Scoring R-CNN [47], Mask R-CNN [12], Hybrid Task Cascade [48] and YOLACT [49]). It also exhibits a competitive performance with models specifically designed for extracting threatening items from X-ray scans (such as CST [25], TST [7], and TSD [50]).

The CIE-Net lags from the fine-tuning-based contour instance segmentation framework TST [7] in terms of  $\mu_{ap}^{b:50}$ . However,

over the incremental TST- $L_s$  [7] version, it achieves 11.29% on the GDXray dataset, 14.65% improvements on the SIXray dataset, and 26.88% on the combined dataset in terms of IoU (see Table 4). The TST [7] possesses the capacity to eliminate unwanted baggage contours due to extensive fine-tuning on the large-scale training datasets, resulting in the better extraction of the threatening items. In return, the TST [7] requires large-scale well-annotated training data to achieve optimal performance. Indeed, when we trained TST [7] framework incrementally on small-scale training batches using the proposed  $L_s$  loss function to compare it with the CIE-Net fairly, it produces degraded performance, as evidenced from the results mentioned above.

Compared to the meta-transfer learning-based baggage threat detector (TSD) [50], our framework achieves 15.62% higher performance in terms of  $\mu_{ap}^{b:50}$  on the SIXray dataset (Table 5). However, on GDXray [44], it lags from the TSD by 6.61%. The superiority of [50] here stems from its capacity to generate the dual-energy tensors [50] that can effectively highlight the transitions of the contraband items from the grayscale X-ray scans. However, TSD is still sensitive to extremely cluttered baggage threats, as evident through its performance on the SIXray [1] dataset.

The performance of CIE-Net, in terms of the  $\mu_{ap}^{b:50}$ , is although lagging from the original CST framework [25] in Table 5. But this comparison is unfair as the original CST [25] framework is non-incremental and uses more training data and computational resources to produce these results. Nevertheless, under fair comparison, the CIE-Net outperforms CST [25] by 4.52% and 3.94% on GDXray [44] and SIXray [1], respectively, in  $\mu_{ap}^{b:50}$  (see Table 5). Also, the CST framework is extremely parametric dependent (i.e., it has to be tuned for each dataset independently). Therefore, it does not generalize well for scans and datasets having drastically varying properties. Furthermore, it also lacks the inherent ability to generate items mask and falls under conventional object detectors.

With regard to run-time performance, the CIE-Net is about two-time faster than several instance segmentation models like MS R-CNN [47], Mask R-CNN [12], and HTC [48]. It also showed a modest performance compared to YOLOv3 [57], CST [25], RetinaNet [2] and YOLCAT [49]. Nonetheless, looking at both accuracy and efficiency figures in, respectively, Table 4, 5, and 6, we can assert that the CIE-Net realizes the best trade-off between time and performance. It is also important to point out that the CIE-Net model's current conception is mainly driven by accurately recognizing the cluttered and overlapping contraband items rather than achieving efficiency. However, we envisage different measures to enhance this aspect in the future. A first remedy can be replacing the conventional convolutional blocks with residual driven atrous convolutions (with variable dilation factors) [58], [59], resulting in a significant reduction of the trainable parameters, thus increasing the overall run-time performance by many folds. Furthermore, we can generate a lightweight version of the CIE-Net by employing a switching mechanism [60] to process only positive regions showcasing contraband items and their instances while ignoring the negative regions. In addition to this, we also envisage employing multi-task

attention networks [30] and adversarial domain adaptation [31] schemes as future work to further improve the threat detection performance of the proposed framework.

## VII. CONCLUSION

This paper presents a novel instance segmentation framework that utilizes incremental learning and a conventional encoder-decoder architecture to extract and recognize heavily cluttered, occluded, and overlapping contraband items from multi-vendor baggage X-ray scans. Since the proposed framework is powered through incremental learning, it reaps the benefit of using small-scale training data and bypasses hectic ground-truth generation mechanisms to make semantic segmentation networks instance-aware. The proposed framework has an in-built capacity to resist catastrophic forgetting through a proposed instance segmentation loss function, introducing a novel feature of incorporating mutual information loss embedding the complex inter-dependencies between old knowledge and newly learned information through Bayesian inference. The proposed scheme is unique as it modifies the conventional semantic segmentation networks to perform instance-aware segmentation via incremental learning. By being trained on two different datasets and their combination, the proposed framework produces the best results compared to existing state-of-the-art solutions in multiple metrics, evidencing the ability to effectively recognize the cluttered and overlapping objects through instance segmentation rather than through object detectors. To the best of our knowledge, it is the only framework to date, which can accurately extract overlapping baggage items from the multi-vendor grayscale and colored X-ray images (in an incremental fashion) despite the significant variations in the scan features of both datasets. In addition to the envisaged task mentioned in the Discussion section to optimize the model design, future work will consider investigating the challenging problem of detecting 3D-printed items (e.g., *guns*). These items, made from organic matter, have low visibility in the X-ray scans. Devising proper models for this category of objects is our potential future work.

## REFERENCES

- [1] C. Miao, L. Xie, F. Wan, C. Su, H. Liu, J. Jiao, and Q. Ye, "SIXray: A Large-scale Security Inspection X-ray Benchmark for Prohibited Item Discovery in Overlapping Images," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2119–2128, 2019.
- [2] T. Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [3] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [4] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, June 2015.
- [5] S. Akçay, M. E. Kundegorski, C. G. Willcocks, and T. P. Breckon, "Using Deep Convolutional Neural Network Architectures for Object Classification and Detection Within X-ray Baggage Security Imagery," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 9, pp. 2203–2215, 2018.
- [6] Y. F. A. Gaus, N. Bhowmik, S. Akçay, P. M. Guillén-García, J. W. Barker, and T. P. Breckon, "Evaluation of a Dual Convolutional Neural Network Architecture for Object-wise Anomaly Detection in Cluttered X-ray Security Imagery," in *2019 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2019.

- [7] T. Hassan and N. Werghe, "Trainable Structure Tensors for Autonomous Baggage Threat Detection Under Extreme Occlusion," *Asian Conference on Computer Vision (ACCV)*, September 2020.
- [8] J. An, H. Zhang, Y. Zhu, and J. Yang, "Semantic Segmentation for Prohibited Items in Baggage Inspection," *International Conference on Intelligence Science and Big Data Engineering. Visual Data Engineering*, pp. 495–505, 2019.
- [9] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon, "GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training," in *Asian Conference on Computer Vision*, pp. 622–637, Springer, 2018.
- [10] E. U. Commission, "List of Prohibited Articles in your Cabin Baggage," Mobility and Transport, 2020.
- [11] Y. F. A. Gaus *et al.*, "Evaluating the Transferability and Adversarial Discrimination of Convolutional Neural Networks for Threat Object Detection and Classification within X-ray Security Imagery," *arXiv preprint arXiv:1911.08966*, 2019.
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *IEEE International Conference on Computer Vision (ICCV)*, pp. 2961–2969, 2017.
- [13] D. Mery, E. Svec, M. Arias, V. Riffo, J. M. Saavedra, and S. Banerjee, "Modern Computer Vision Techniques for X-Ray Testing in Baggage Inspection," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, Volume: 47, Issue: 4, pp. 682–692 April 2017.
- [14] S. Akçay and T. Breckon, "Towards Automatic Threat Detection: A Survey of Advances of Deep Learning within X-ray Security Imaging," *preprint arXiv:2001.01293*, 2020.
- [15] M. Bastan, W. Byeon, and T. Breuel, "Object Recognition in Multi-View Dual Energy X-ray Images," *British Machine Vision Conference*, 2013.
- [16] M. Bastan, "Multi-view Object Detection In Dual-energy X-ray Images," *Machine Vision and Applications*, p. 1045–1060, 2015.
- [17] G. Heitz and G. Chechik, "Object Separation in X-ray Image Sets," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, p. 2093–2100, 2010.
- [18] M. Bastan, M. R. Yousefi, and T. M. Breuel, "Visual Words on Baggage X-ray Images," 14th International Conference on Computer Analysis of Images and Patterns, August 2011.
- [19] M. E. Kundegorski, S. Akçay, M. Devereux, A. Mouton, and T. P. Breckon, "On using Feature Descriptors as Visual Words for Object Detection within X-ray Baggage Security Screening," in *IEEE International Conference on Imaging for Crime Detection and Prevention (ICDP)*, 2016.
- [20] D. Mery, E. Svec, and M. Arias, "Object Recognition in Baggage Inspection Using Adaptive Sparse Representations of X-ray Images," in *Pacific-Rim Symposium on Image and Video Technology*, p. 709–720, 2016.
- [21] V. Riffo and D. Mery, "Automated Detection of Threat Objects Using Adapted Implicit Shape Model," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, Volume: 46, Issue: 4, pp. 472–482, 2016.
- [22] J. Redmon and A. Farhadi, "YOLO9000: Better, Faster, Stronger," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [23] N. Jaccard, T. W. Rogers, E. Morton, and L. D. Griffin, "Detection of Concealed Cars in Complex Cargo X-ray Imagery using Deep Learning," in *Journal of X-Ray Science and Technology*, p. 323–339, 2017.
- [24] L. Zou, T. Yusuke, and I. Hitoshi, "Dangerous Objects Detection of X-ray Images Using Convolution Neural Network," *Security with Intelligent Computing and Big-data Services*, 2018.
- [25] T. Hassan, M. Bettayeb, S. Akçay, S. Khan, M. Bennamoun, and N. Werghe, "Detecting Prohibited Items in X-ray Images: A Contour Proposal Learning Approach," *IEEE International Conference on Image Processing (ICIP)*, pp. 2016–2020, 2020.
- [26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [27] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation," *European Conference on Computer Vision (ECCV)*, 2018.
- [28] T. Akilan, Q. M. J. Wu, and W. Zhang, "Video foreground extraction using multi-view receptive field and encoder–decoder dcnn for traffic and surveillance applications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 10, pp. 9478–9493, 2019.
- [29] L. Wang, D. Li, Y. Zhu, L. Tian, and Y. Shan, "Dual super-resolution learning for semantic segmentation," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3773–3782, 2020.
- [30] Q. Wang, T. Han, Z. Qin, J. Gao, and X. Li, "Multitask Attention Network for Lane Detection and Fitting," *IEEE Transactions on Neural Networks and Learning Systems*, December 2020.
- [31] Q. Wang, J. Gao, and X. Li, "Weakly Supervised Adversarial Domain Adaptation for Semantic Segmentation in Urban Scenes," *IEEE Transactions on Image Processing*, 2019.
- [32] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon, "Skip-GANomaly: Skip Connected and Adversarially Trained Encoder-Decoder Anomaly Detection," *arXiv preprint arXiv:1901.08954*, 2019.
- [33] G. Hinton, O. Vinyals, and J. Dean, "Distilling the Knowledge in a Neural Network," *arXiv:1503.02531*, March 2015.
- [34] R. Aljundi, P. Chakravarty, and T. Tuytelaars, "Expert Gate: Lifelong Learning with a Network of Experts," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [35] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental Classifier and Representation Learning," *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [36] Y. Tian, D. Krishnan, and P. Isola, "Contrastive Representation Distillation," *International Conference on Learning Representations*, 2020.
- [37] J. H. Cho and B. Hariharan, "On the Efficacy of Knowledge Distillation," *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [38] D. Lopez-Paz and M. Ranzato, "Gradient Episodic Memory for Continual Learning," *Neural Information Processing Systems*, 2017.
- [39] U. Michieli and P. Zanuttigh, "Knowledge Distillation for Incremental Learning in Semantic Segmentation," *arXiv:1911.03462*, 2020.
- [40] S. Zagoruyko, A. Lerer, T.-Y. Lin, P. O. Pinheiro, S. Gross, S. Chintala, and P. Dollár, "A MultiPath Network for Object Detection," *arXiv:1604.02135v2*, August 2016.
- [41] D. Haris and S. Harris, "Digital design and computer architecture." *San Francisco, Calif.: Morgan Kaufmann*, p. 129. ISBN 978-0-12-394424-5, August 2012.
- [42] Z. Chen and B. Liu, "Continual Learning and Catastrophic Forgetting," In "Lifelong Machine Learning", *Morgan & Claypool Publishers*, 2018.
- [43] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid Scene Parsing Network," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2881–2890, 2017.
- [44] D. Mery, V. Riffo, U. Zscherpel, G. Mondragón, I. Lillo, I. Zuccar, H. Lobel, and M. Carrasco, "GDxray: The database of X-ray images for nondestructive testing," *Journal of Nondestructive Evaluation*, Volume 34, Issue: 4, 2015.
- [45] M. D. Zeiler, "ADADELTA: An Adaptive Learning Rate Method," *arXiv:1212.5701*, 2012.
- [46] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *European Conference on Computer Vision (ECCV)*, 2014.
- [47] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang, "Mask Scoring R-CNN," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6409–6418, 2019.
- [48] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang, *et al.*, "Hybrid Task Cascade for Instance Segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4974–4983, 2019.
- [49] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "YOLACT: Real-Time Instance Segmentation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 9157–9166, 2019.
- [50] T. Hassan, M. Shafay, S. Akçay, S. Khan, M. Bennamoun, E. Damiani, and N. Werghe, "Meta-Transfer Learning Driven Tensor-Shot Detector for the Autonomous Localization and Recognition of Concealed Baggage Threats," *MDPI Sensors*, November 2020.
- [51] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [52] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," *arXiv:1505.04597*, 2015.
- [53] J. Long, E. Shelhamer, and T. Darrell, "Fully Convolutional Networks for Semantic Segmentation," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3431–3440, 2015.
- [54] U. Michieli and P. Zanuttigh, "Incremental Learning Techniques for Semantic Segmentation," *IEEE International Conference on Computer Vision Workshops (ICCVW)*, 2019.
- [55] J. Zhang, J. Zhang, S. Ghosh, D. Li, S. Tasci, L. Heck, H. Zhang, and C.-C. J. Kuo, "Class-incremental Learning via Deep Model Consolidation," *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020.
- [56] F. Tung and G. Mori, "Similarity-Preserving Knowledge Distillation," *IEEE International Conference on Computer Vision (ICCV)*, 2019.

- [57] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," *arXiv*, 2018.
- [58] H. Raja, T. Hassan, M. U. Akram, and N. Werghe, "Clinically Verified Hybrid Deep Learning System for Retinal Ganglion Cells Aware Grading of Glaucomatous Progression," *IEEE Transactions on Biomedical Engineering*, 2020.
- [59] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell, "Understanding Convolution for Semantic Segmentation," IEEE Winter Conference on Applications of Computer Vision (WACV), 2018.
- [60] H. Chen, H. Lin, and M. Yao, "Improving the Efficiency of Encoder-Decoder Architecture for Pixel-Level Crack Detection," IEEE Access, December 2019.