

ZBS: Zero-shot Background Subtraction via Instance-level Background Modeling and Foreground Selection

Yongqi An^{1,2} Xu Zhao^{1,*} Tao Yu^{1,2} Haiyun Guo^{1,2}
 Chaoyang Zhao¹ Ming Tang^{1,2} Jinqiao Wang^{1,2}

National Laboratory of Pattern Recognition, Institute of Automation, CAS, Beijing, China¹
 School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China²

{yongqi.an, xu.zhao, haiyun.guo, tangm, jqwang}@nlpr.ia.ac.cn
 yutao2022@ia.ac.cn

Abstract

Background subtraction (BGS) aims to extract all moving objects in the video frames to obtain binary foreground segmentation masks. Deep learning has been widely used in this field. Compared with supervised-based BGS methods, unsupervised methods have better generalization. However, previous unsupervised deep learning BGS algorithms perform poorly in sophisticated scenarios such as shadows or night lights, and they cannot detect objects outside the pre-defined categories. In this work, we propose an unsupervised BGS algorithm based on zero-shot object detection called Zero-shot Background Subtraction (ZBS). The proposed method fully utilizes the advantages of zero-shot object detection to build the open-vocabulary instance-level background model. Based on it, the foreground can be effectively extracted by comparing the detection results of new frames with the background model. ZBS performs well for sophisticated scenarios, and it has rich and extensible categories. Furthermore, our method can easily generalize to other tasks, such as abandoned object detection in unseen environments. We experimentally show that ZBS surpasses state-of-the-art unsupervised BGS methods by 4.70% F-Measure on the CDnet 2014 dataset. The code is released at <https://github.com/CASIA-IVA-Lab/ZBS>.

1. Introduction

Background subtraction (BGS) is a fundamental task in computer vision applications [7], such as autonomous navigation, visual surveillance, human activity recognition, etc [15]. BGS aims to extract all moving objects as foreground in each video frame and outputs binary segmenta-

*Corresponding Author

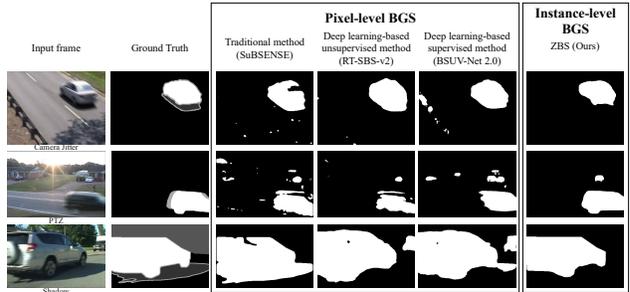


Figure 1. The performance of different BGS methods. Previous BGS methods based on pixel-level background models may misjudge noisy background as foreground objects, such as camera-jitter, PTZ, and shadow. Our method based on an instance-level background model can obtain precise foreground edges, effectively reducing the confusion of background pixels as foreground objects.

tions.

The most straightforward BGS algorithm is to directly compare the current frame with the "stationary" background image [7]. However, this strategy cannot handle complex scenarios, such as dynamic background, illumination changes, and shadows. Therefore, more sophisticated BGS techniques [7, 20, 24, 48] have been proposed in the past decades. The traditional methods improve performance in two aspects. The first is to design more robust feature representations, including color features [44], edge features [20], motion features [48], and texture features [12]. The second is to design more suitable background models, such as Gaussian mixture models [36], kernel density estimation models [14], CodeBook [21], ViBe [4], SuBSENSE [34], and PAWCS [35]. The traditional methods have relatively adequate generalization capacity since they are not optimized on specific scenarios or categories of objects. How-

ever, these methods only utilize hand-craft features to determine whether each pixel belongs to the foreground. We call these methods pixel-level BGS since they use pixel-based or local pixels-based background models. They are sensitive to natural variations such as lighting and weather.

Over the years, deep learning-based BGS algorithms have been proposed, including supervised BGS and unsupervised BGS. Supervised BGS algorithms have achieved satisfactory performance on CDnet 2014 benchmark [11, 24, 31, 41, 46]. However, these methods usually have to be trained on the first several frames of the test videos, which limits the application to unseen scenarios. Unsupervised algorithms overcome this shortcoming. Most of them combine semantic segmentation models into traditional BGS algorithms. These algorithms pre-select 12 categories as foreground from 150 categories of semantic segmentation models [9]. Existing state-of-the-art unsupervised methods still detect night light and heavy shadows as foreground objects. As shown in Figure 1, it is difficult for pixel-level background model to accurately distinguish the edges of foreground objects.

To tackle the above problems, we propose a novel background subtraction framework based on zero-shot object detection (ZBS). The zero-shot object detection, or also named open-vocabulary object detection, aims to detect unseen objects outside of the pre-defined categories [49]. Figure 2 shows the framework of our method. The method includes all-instance detection, instance-level background modeling, and foreground instance selection. In the all-instance detection stage, any zero-shot detector can be used. We use a zero-shot object detection model named Detic [49] as the all-instance detector to transform the raw image pixels into structured instance representations, including categories, boxes, and masks. In the background modeling stage, our method builds an instance-level background model based on the motion information of instances. If an object is stationary, our algorithm adds it to the background model. In the foreground instance selection stage, the proposed algorithm selects the output of the all-instance detector when the new frame comes. If the instance complies with Rule 2 in Figure 2 (c), it is the foreground in the final binary mask. Benefiting from the full use of instance information, our instance-level BGS method performs better in complex scenarios, such as shadows, camera jitter, night scenes, *etc.* ZBS rarely detects noisy background as foreground objects by mistake. Due to the characteristics of the detector, the proposed method can detect most of the categories in the real world and can detect the unseen foreground categories outside the pre-defined categories. ZBS achieves remarkably 4.70% F-Measure improvements over state-of-the-art unsupervised methods.

Our main contributions are listed as follows:

- We propose a novel background subtraction frame-

work that has the instance-level background model;

- The proposed framework uses a zero-shot object detection model to obtain a more general and generalized deep learning-based unsupervised BGS algorithm;
- Our method achieves the state-of-the-art in all unsupervised BGS methods on the CDnet 2014 dataset.

2. Related work

2.1. Deep learning-based Supervised Methods

Deep learning methods have been widely used for BGS due to their ability to learn high-level representations from training data [8]. Braham *et al.* [10] presented the first work using deep learning for background subtraction. FgSegNet [24] is a representative work that focuses on learning multi-scale features for foreground segmentation. CascadeCNN [41] employs a cascade structure to synthesize the basic CNN model and the multi-scale CNN model. Zhao *et al.* [46] propose an end-to-end two-stage deep CNN to reconstruct the background and separate the foreground from the background jointly. Chen *et al.* [11] and Sakkos *et al.* [31] use ConvLSTM and 3DCNN, respectively, to process spatio-temporal information. In addition, Siamese neural networks [18, 32], generative adversarial networks (GAN) [1, 2, 47], and autoencoders (AE) [33] have also been employed for BGS.

Recently, [37, 38, 42, 45] demonstrated better generality for unseen videos with training on limited data. However, these models are trained only on datasets containing a few categories and scenes, limiting their ability to cope with more complex real-world detection and segmentation tasks.

2.2. Semantic background subtraction

SemanticBGS [9] is the first motion detection framework to utilize object-level semantics for improving background subtraction. By combining semantic segmentation and background subtraction algorithms, it significantly reduces false positive detections and effectively identifies camouflaged foreground objects. RTSS [43] performs foreground detection and semantic segmentation in a parallel manner, using the semantic probability of pixels to guide the construction and update of the background model. This method achieves real-time semantic background subtraction. RT-SBS [13] adopts a similar approach and improves the performance, achieving a real-time semantic background subtraction algorithm at 25 frames per second.

Despite their advancements, semantic background subtraction methods are still fundamentally pixel-level background models. All of these semantic BGS methods necessitate a predefined list of foreground classes, which require expert knowledge and pose challenges for implementation in various scenarios. Furthermore, the limited number of categories in semantic segmentation networks (up to

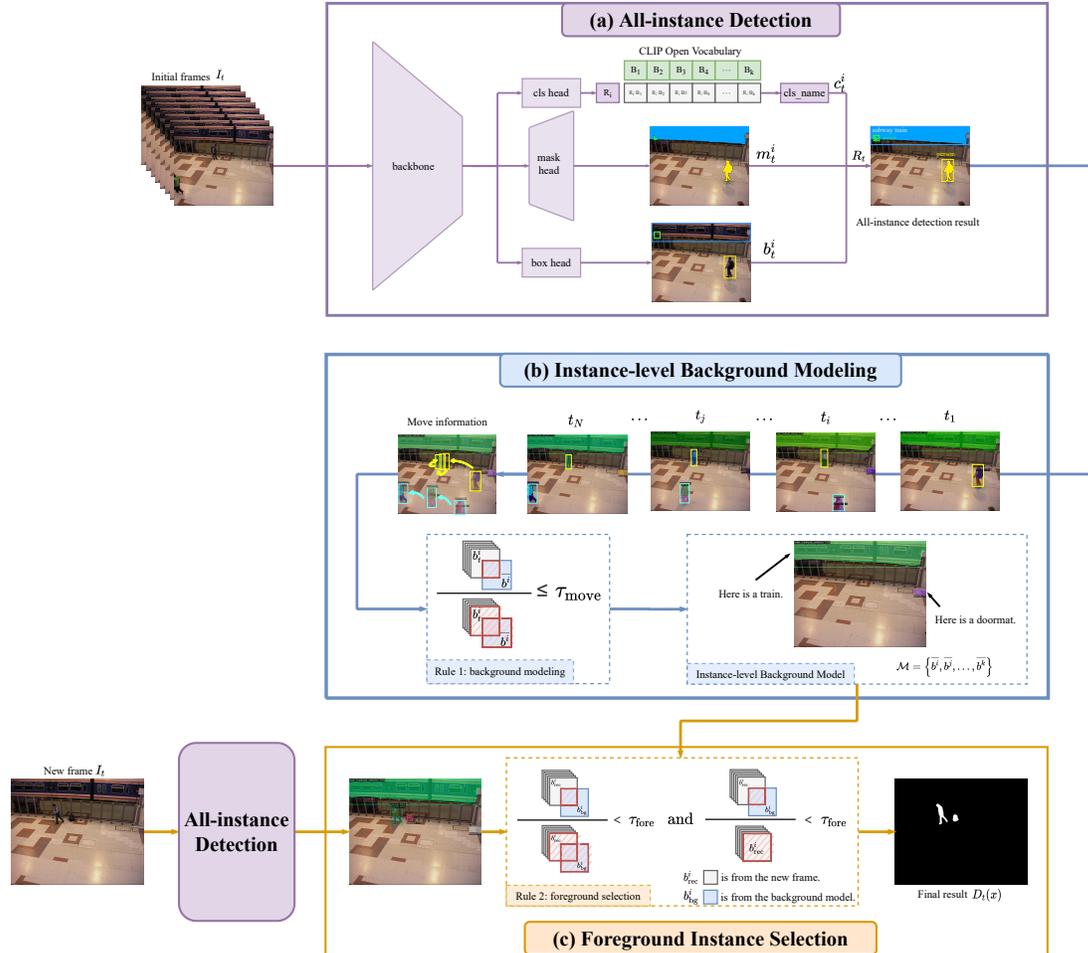


Figure 2. The framework of ZBS. (a) All-instance detection. We use a zero-shot object detection model named Detic [49] to transform the pixel-level image into a structured-instance representation, including categories, boxes, and masks. Specifically, the categories are obtained by CLIP. (b) Instance-level background modeling. The proposed method analyzes the motion information of instances. If the instance complies with Rule 1, the stationary instance will be added to the background model. (c) The new frame output by Detic will be compared with the instance-level background model. If the instance complies with Rule 2, it will be the foreground in the final result.

150 categories) hinders their ability to detect moving foregrounds in an open-vocabulary setting, an aspect that is becoming increasingly important in today’s environment.

The proposed ZBS builds an instance-level background model capable of detecting most real-world categories without the need for predefined foreground classes, thus offering greater practicality.

2.3. Zero-shot object detection

In the era of deep learning, supervised learning has achieved outstanding performance on many tasks, but too much training data has to be used. Moreover, these methods can not classify or detect the objects of categories outside the training datasets’ annotations. To solve this problem, zero-shot learning [3] was proposed, hoping to classify images of categories that were never seen in the train-

ing process. The zero-shot object detection developed from this aims to detect objects outside the training vocabulary. Earlier work studied exploiting attributes to encode categories as vectors of binary attributes and learning label embeddings [3]. The primary solution in Deep Learning is to replace the last classification layer with a language embedding of class names (e.g., GloVe [27]). Rahman *et al.* [29] and Li *et al.* [23] improve by introducing external textual information Classifier Embeddings. ViLD [17] upgrades language embeddings to CLIP [28] and extracts regional features from CLIP image features. Detic [49] also adopts CLIP as a classifier and uses additional image-level data for joint training, dramatically expanding the number of categories and performance of detection-level tasks. This paper uses the Detic detector for the BGS algorithm.

3. Method

3.1. Overview

ZBS is among the novel unsupervised BGS algorithms for real applications. It is a zero-shot object detection based on the model that can obtain an instance-level background model. ZBS contains three stages: all-instance detection, instance-level background modeling, and foreground instance selection. Figure 2 illustrates the framework of our method. First, ZBS uses an all-instance detection model to acquire the structured-instance representation. Then, an instance-level background model is built and maintained through the movement information of instances. Finally, when a new frame comes, we will select the moving foreground from the detector outputs based on the background model. We convert the result into a binary mask to compare with other BGS methods.

Algorithm 1 : The ZBS algorithm process.

- 1: Initialize the zero-shot detector as \mathcal{Z}
 - 2: Initialize the background model as \mathcal{M}
 - 3: **while** current frame I_t is valid **do**
 - 4: **Stage 1: All-instance detection**
 - 5: output the result $R_t \leftarrow \mathcal{Z}(I_t)$
 - 6: **Stage 2: Instance-level background model**
 - 7: get the track of each instance from b_t (part of R_t)
 - 8: calculate the IoU_{\min} of $b_0^i \dots b_t^i$ and \bar{b}^i
 - 9: update \mathcal{M} based on IoU_{\min} and τ_{move}
 - 10: **Stage 3: Foreground instance selection**
 - 11: separate \mathcal{M} and b_t by instance-id
 - 12: calculate the IoU and IoF of \mathcal{M} and b_t^i
 - 13: get a binary mask $D_t(x)$ based on IoU&IoF and τ_{fore}
 - 14: current frame \leftarrow next frame
- end**
-

3.2. All-instance Detection

The goal of background subtraction is to extract all moving objects as foreground in each video frame. Traditional unsupervised BGS methods rely on pixel-level background models, which struggle to differentiate noisy backgrounds from foreground objects. To address this, we propose an instance-level background model. It utilizes an instance detector to locate the objects of all possible categories and all locations in the image and convert the raw image pixels into structured instance representation.

Intuitively, most existing trained instance segmentation networks can be used. Besides, the categories of the training datasets adapt to most domain-adapted object detection scenarios. However, instance segmentation networks cannot detect and segment the objects of categories outside the training datasets’ annotations.

Recently, with the development of self-supervised training and the foundation models [28], several practical zero-shot object detection methods have been proposed [17, 49]. These methods can detect almost thousands of categories of objects without being trained on the applied scenarios. Therefore, to obtain a more general and generalized deep learning background modeling method, we adopt the zero-shot object detection method Detic [49] as the detector of our BGS method. Detic [49] can detect 21k categories of objects and segment the object’s masks.

Distinguished from the instance segmentation, we call this process *all-instance detection*. After the all-instance detection stage, the video frame I_t is structured by zero-shot detector \mathcal{Z} as instance representation R_t . The representation R_t includes instance boxes b_t^i and segmentation masks m_t^i with category labels c_t^i , where i is the *id* of the instance, and t refers to the t -th frame of the video.

3.3. Instance-level Background Modeling

Based on the all-instance detection results, the algorithm should distinguish which objects have moved and which have not. The ideal instance-level background model should be the collection of all stationary object instances. It is the basis for foreground instance selection. We define the background model \mathcal{M} as:

$$\mathcal{M} = \left\{ \bar{b}^i, \bar{b}^j, \dots, \bar{b}^k \right\} \quad (1)$$

The instance-level background model \mathcal{M} is a collection of detection boxes for static instances, $\bar{b}^i, \bar{b}^j, \bar{b}^k$ are static instances with $id = i, j, k$, and the value is the average of the coordinates of all boxes in the past trajectory of this instance. It reflects which locations have stationary instances and which are background without instances. As shown in Figure 2 (b), our method uses the initial frames to obtain an initial instance-level background model and update the background model with a certain period in subsequent frames ($\Delta T = 100$ is chosen in this paper).

The details are shown in Algorithm 1. There are three steps for the instance-level background modeling stage. First, the proposed method utilizes the detector \mathcal{Z} output boxes b_t^i from past frames to obtain the tracks of each instance (tracks are obtained by SORT method [5]). Second, ZBS computes the average value of the coordinates for the upper-left and lower-right corners of each bounding box within the corresponding trajectory of the instance, denoted as \bar{b}^i . Then we can obtain the minimum value of IoU of b_t^i and \bar{b}^i , which means the maximum movement between the positions compared to the average in the whole trajectory. In our implementation, we apply a median filter to the trajectory IoU. This helps mitigate abrupt changes in IoU caused by object occlusion. Experiments in Table 2 demonstrate that this improvement is beneficial. Finally, the

update strategy of the instance-level background model is as Equation (2):

$$\mathcal{M} = \begin{cases} \mathcal{M} \cup \bar{b}^i, & \text{if } \text{IoU}_{\min}(b_t^i, \bar{b}^i) \geq \tau_{\text{move}} \\ \mathcal{M} - (\mathcal{M} \cap \bar{b}^i), & \text{otherwise.} \end{cases} \quad (2)$$

where b_t^i denotes the i -th instance in t -th image frame. \bar{b}^i denotes the average of all boxes \bar{b}^i for each instance. τ_{move} is the threshold for judging whether the instance is moving. If it remains stationary, put it into the background model \mathcal{M} ; otherwise, remove it from the background model \mathcal{M} .

Implementation Details. To build a more robust background model, we choose a smaller τ_{conf} ¹ in the instance-level background modeling stage, which helps more stationary instances to be incorporated into the background model \mathcal{M} , it is called Δconf .

3.4. Foreground Instance Selection

The key to foreground instance selection is accurately judging whether the instance has moved compared to the background model. Object occlusion is a common challenge. When objects are occluded, the object behind them can easily be misjudged as moving. To balance sensitivity and robustness to object occlusion, we introduce IoF (Intersection over Foreground) as a complement to IoU (Intersection over Union), which is calculated as Equation (3):

$$\text{IoF} = \frac{b_{\text{rec}}^i \cap b_{\text{bg}}^i}{b_{\text{rec}}^i}. \quad (3)$$

where b_{rec}^i denotes the instance in the recent frame, b_{bg}^i denotes the i -th instance in the instance-level background model \mathcal{M} .

Figure 3 shows how IoF works. If the instance is not moving but is obscured, IoU drops considerably while IoF still preserves a high value. By judging the two metrics together, it is robust to determine whether a certain instance is a foreground.

As shown in the foreground instance selection stage of Algorithm 1, the detection result of each current frame is matched with the set of object instances contained in the instance-level background model. The proposed method uses two metrics, IoU and IoF, to determine whether the instance can be used as a foreground. If both IoU and IoF are smaller than the foreground selection threshold τ_{fore} , the instance is new or has moved and should be considered as foreground. On the contrary, if either IoU or IoF is larger than the threshold τ_{fore} , the instance should not be considered as foreground. The rule of foreground instance selec-

¹ τ_{conf} is a score threshold from [49] for the outputs of the all-instance detection stage. This threshold determines the confidence level of the outputs in the first stage.

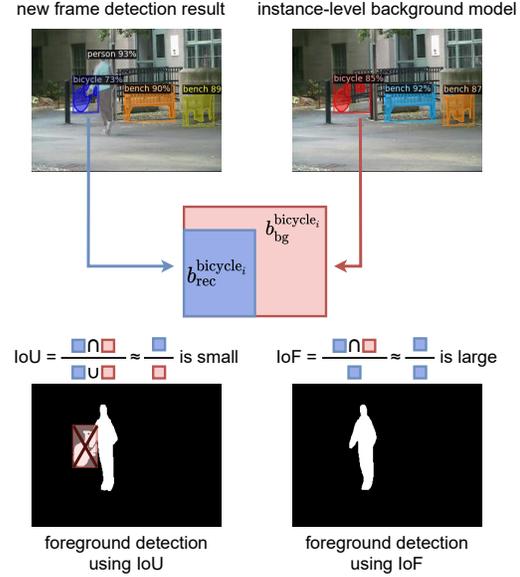


Figure 3. A typical case for object obstruction. A person half obscures the bicycle on the left side of the image. The IoU of $b_{\text{rec}}^{\text{bicycle}_i}$ and $b_{\text{bg}}^{\text{bicycle}_i}$ is very small, while the IoF of $b_{\text{rec}}^{\text{bicycle}_i}$ and $b_{\text{bg}}^{\text{bicycle}_i}$ is still large.

tion can be expressed by Equation (4).

$$D_t(x) = \begin{cases} FG, & \text{if } \text{IoU}(b_{\text{rec}}^i, b_{\text{bg}}^i) < \tau_{\text{fore}} \\ & \text{and } \text{IoF}(b_{\text{rec}}^i, b_{\text{bg}}^i) < \tau_{\text{fore}}; \\ BG, & \text{otherwise.} \end{cases} \quad (4)$$

where $D_t(x)$ is regarded as the x -th instance in t -th frame whether should be a foreground, FG means foreground, BG means background. b_{rec}^i denotes the box of i -th instance in recent frame. b_{bg}^i denotes the box of foreground instance in the instance-level background model \mathcal{M} .

4. Experiments

4.1. Dataset and Evaluation Metrics

We evaluate the performance of the proposed method on the CDnet 2014 dataset [40]. The CDnet 2014 dataset is the most famous benchmark for change detection, including 53 video sequences and 11 categories corresponding to different classic application scenarios. The main categories of the dataset include *Bad Weather*, *Baseline*, *Camera Jitter*, *Dynamic Background*, *Intermittent Object Motion*, *Low Framerate*, *Night videos*, *Pan-Tilt-Zoom*, *Shadow*, *Thermal*, and *Turbulence*. Ground Truth annotated manually is available for every frame in video sequences and tells us whether each pixel belongs to the background or the foreground. Specifically, in the Ground Truth, Static, Shadow, Non-Region of Interest (Non-ROI), Unknown, and Moving pixels are assigned respectively to grayscale values 0, 50, 85, 170, and 255. We select Recall (Re), Precision (Pr) and

F-Measure ($F-M$) to evaluate the performance on the CD-net 2014 dataset.

Following [16], we regard Static and Shadow pixels as negative samples (background), regard Moving pixels as positive samples (foreground), and ignore Non-ROI and Unknown pixels to ensure the fairness of the metrics.

4.2. Hyper-parameter Sensitivity Analysis

As mentioned earlier, our method requires setting several parameters before use: the threshold for all-instance detection τ_{conf} , the threshold for background modeling τ_{move} , and the threshold for foreground selection τ_{fore} .

Different parameters suit different situations. The large τ_{conf} is better for more straightforward scenarios. The large τ_{move} is better for fast motion or low frame rate scenarios. The small τ_{fore} is robust for camera jitter scenarios. The performance with different parameter settings is shown in Figure 4. For most scenarios, the sensitivity of the two parameters τ_{move} , τ_{fore} is low, and the impact of the changes of these two hyper-parameters on F-Measure fluctuates within $\pm 1\%$. When τ_{fore} is equal to 1, the foreground instance selector treats all new instances as foreground, so the precision and F-Measure have a big drop. τ_{conf} determines the output of the zero-shot detector, which is very important for the subsequent two stages, and different scenes often apply to different thresholds. The universal parameters are $\tau_{\text{conf}} = 0.6$, $\tau_{\text{move}} = 0.5$, $\tau_{\text{fore}} = 0.8$ in Experiments.

4.3. Quantitative Results

Table 1 shows the comparison of our method among other BGS algorithms, in which F-Measure is observed. These algorithms could be classified into two parts: supervised and unsupervised algorithms. Most supervised algorithms, such as FgSegNet [24] and CascadeCNN [41], have nearly perfect F-Measure because they are trained with some frames in test videos. However, the performance of these methods decreases significantly when applied to unseen videos because of the lack of generalization ability. FgSegNet in unseen videos only achieves 0.3715 F-Measure. STPNet [42] and BSUV-Net 2.0 [37] are supervised algorithms designed explicitly for unseen videos and can achieve F-Measure of around 0.8. IUTIS-5 [6] is a special supervised algorithm that learns how to combine various unsupervised algorithms from datasets.

The remaining methods are unsupervised algorithms [19, 22, 34, 35, 43] which naturally can handle unseen videos. The results show that our method outperforms all unsupervised algorithms. In particular, ZBS outperforms the state-of-the-art RT-SBS-v2 [13] by 4.70%. Moreover, ZBS outperforms supervised method in unseen videos BSUV-Net 2.0 [37]. When considering per-category F-Measure, our method has advantages in seven out of eleven categories, such as *Camera Jitter*, *Intermittent Object Motion*,

and *Night Videos*. However, our method cannot deal with *Turbulence* well because the detector of the all-instance detection module cannot adapt to the unnatural image distribution of *Turbulence* scenarios without training.

4.4. Ablation Study

Ablation experiments are conducted, in which we add the ablation components one by one to measure their effectiveness. The results are summarized in Table 2 with precision, recall, and F-Measure.

The baseline is to use the result of the all-instance detector directly as the foreground. In the instance-level background modeling stage, we only build an instance-level background model, but do not judge whether the foreground is moving. The performance of the algorithm is slightly improved and exceeds the baseline. In the foreground selection stage, the algorithm uses the background model to determine whether the foreground is moving. The performance is greatly improved. Moreover, we propose three modules to enhance the instance-level background model. After adding them to the algorithm one by one, the algorithm’s performance larger gains and the advantages of the instance-level background model are fully demonstrated.

Table 2 shows that simple instance-level background and foreground capture most potential moving objects. The former exhibits higher recall but slightly lower precision than pixel-level background and foreground. Δconf enhances overall performance. Object occlusion impacts the background model and complicates foreground selection. This issue is unique to instance-level representation. We propose the "IoU filter" and "IoF" to mitigate this problem, both of which reduce false positives, particularly "IoF".

4.5. More Analysis

Visual Result. A visual comparison from different methods is shown in Figure 5. It includes six challenging scenarios from the CDnet 2014 dataset. In the "highway" scenario, the main challenge is the shadow of the car and the tree branches moving with the wind. Because of the instance-level foreground detection, our method is robust to noisy background regions. In the "boulevard" scenario, affected by the jitter of the camera, many false positives are produced by other BGS methods. However, our method is robust to camera shake due to the instance-level background model. In the "boats" scenario, GMM and SuBSENSE produce many false positives because of water rippling. ZBS can extract clear detection results within the background disturbance. In the "sofa" scenario, which contains intermittent object motion, the proposed method has better segmentation results at the contour regions of objects. In the "peopleInShade" scenario, ZBS excels in shadow regions. Despite challenges in the "continuousPan" scenario, our method remains robust. These results highlight the advan-

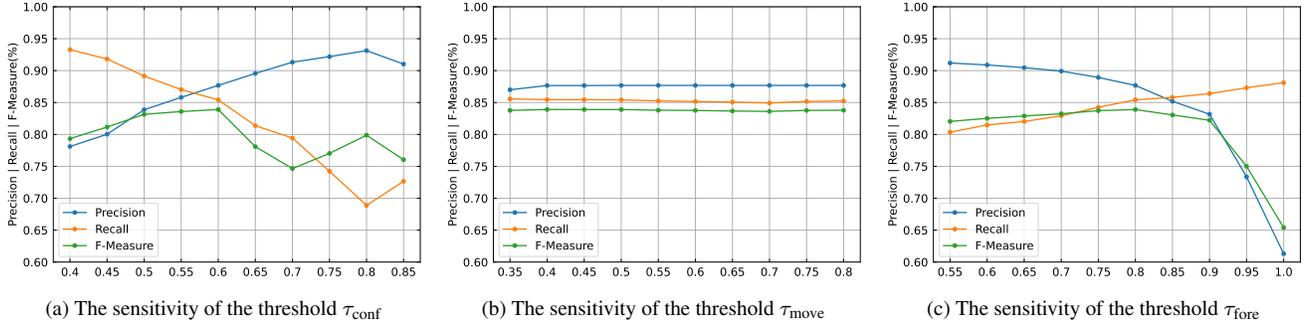


Figure 4. The hyper-parameter sensitivity analysis. The relationship between the thresholds (τ_{conf} , τ_{move} , τ_{fore}) and the evaluation metrics (Precision, Recall, F-Measure) of the BGS algorithms.

Table 1. Overall and per-category F-Measure comparison of different BGS methods on the CDnet 2014 dataset.

Method	baseline	camjitt	dynbg	intmot	shadow	thermal	badwea	lowfr	night	PTZ	turbul	Overall
Supervised algorithms												
CascadeCNN [41]	0.9786	0.9758	0.9658	0.8505	0.9593	0.8958	0.9431	0.8370	0.8965	0.9168	0.9108	0.9209
MU-Net2 [30]	0.9900	0.9824	0.9892	0.9894	0.9845	0.9842	0.9343	0.8706	0.8362	0.8185	0.9272	0.9369
BSPVGAN [47]	0.9837	0.9893	0.9849	0.9366	0.9849	0.9764	0.9644	0.8508	0.9001	0.9486	0.9310	0.9501
FgSegNetv2 [25]	0.9978	0.9971	0.9951	0.9961	0.9955	0.9938	0.9904	0.9336	0.9739	0.9862	0.9727	0.9847
FgSegNet [24] (unseen video)	0.6926	0.4266	0.3634	0.2002	0.5295	0.6038	0.3277	0.2482	0.2800	0.3503	0.0643	0.3715
STPNet [42]	0.9587	0.7721	0.8058	0.8267	0.9114	0.8688	0.8898	0.7297	0.6961	0.6076	0.7248	0.7992
BSUV-Net 2.0 [37]	0.9620	0.9004	0.9057	0.8263	0.9562	0.8932	0.8844	0.7902	0.5857	0.7037	0.8174	0.8387
IUTIS-5 [6]	0.9567	0.8332	0.8902	0.7296	0.8766	0.8303	0.8248	0.7743	0.5290	0.4282	0.7836	0.7717
Unsupervised algorithms												
PAWCS [35]	0.9397	0.8137	0.8938	0.7764	0.8913	0.8324	0.8152	0.6588	0.4152	0.4615	0.6450	0.7403
SuBSENSE [34]	0.9503	0.8152	0.8177	0.6569	0.8986	0.8171	0.8619	0.6445	0.5599	0.3476	0.7792	0.7408
WisenetMD [22]	0.9487	0.8228	0.8376	0.7264	0.8984	0.8152	0.8616	0.6404	0.5701	0.3367	0.8304	0.7535
SWCD [19]	0.9214	0.7411	0.8645	0.7092	0.8779	0.8581	0.8233	0.7374	0.5807	0.4545	0.7735	0.7583
SemanticBGS [9]	0.9604	0.8388	0.9489	0.7878	0.9478	0.8219	0.8260	0.7888	0.5014	0.5673	0.6921	0.7892
RTSS [43]	0.9597	0.8396	0.9325	0.7864	0.9551	0.8510	0.8662	0.6771	0.5295	0.5489	0.7630	0.7917
RT-SBS-v2 [13]	0.9535	0.8233	0.9217	0.8946	0.9497	0.8697	0.8279	0.7341	0.5629	0.5808	0.7315	0.8045
ZBS (Ours)	0.9653	0.9545	0.9290	0.8758	0.9765	0.8698	0.9229	0.7433	0.6800	0.8133	0.6358	0.8515

Table 2. Ablation of the three stages of our methods and other improvements. AID: All-instance detection (Section 3.2). IBM: Instance-level background modeling (Section 3.3). FIS: Foreground instance selection (Section 3.4). Δconf : Different confidence thresholds for background modeling and foreground selection. Filter: Median filtering on movement information. IoF: Intersection over Foreground measurement standard.

AID	IBM	FIS	Enhance		IoF	Pr	Re	F-M
			Δconf	Filter				
✓						0.4076	0.8869	0.4980
✓	✓					0.5343	0.8022	0.5752
✓	✓	✓				0.7468	0.7625	0.7152
✓	✓	✓	✓			0.7529	0.7851	0.7415
✓	✓	✓	✓	✓		0.8249	0.7829	0.7836
✓	✓	✓	✓	✓	✓	0.8802	0.8403	0.8515

tages of our instance-level background-based ZBS.

Runtime Efficiency. Achieving real-time performance is vital for BGS algorithms. The main time-consuming aspect of ZBS is concentrated in the first stage, which involves pre-

training the zero-shot object detection model. We have implemented the two subsequent stages in C++. The FPS is about 20 on one A100 GPU. In the all-instance detection stage, we used parallel computation with a batch size of 8 and processed results sequentially in the subsequent stages. Ultimately, we achieve approximately 44 FPS on an A100 GPU.

Moreover, adopting TensorRT SDK [39], quantization, and frame skipping can further improve the FPS in real-world applications. However, this paper mainly focus on enhancing the accuracy of BGS (F-Measure). In future studies, we plan to further improve its runtime efficiency.

Performance in complex scenarios. To further demonstrate the good performance of ZBS in complex scenarios, we compare the nbShadowError and FPR-S of different BGS methods in the *Shadow* category. FPR-S is calculated as Equation (5). Table 3 shows that our method has an extremely low false positive rate on shadow pixels. ZBS

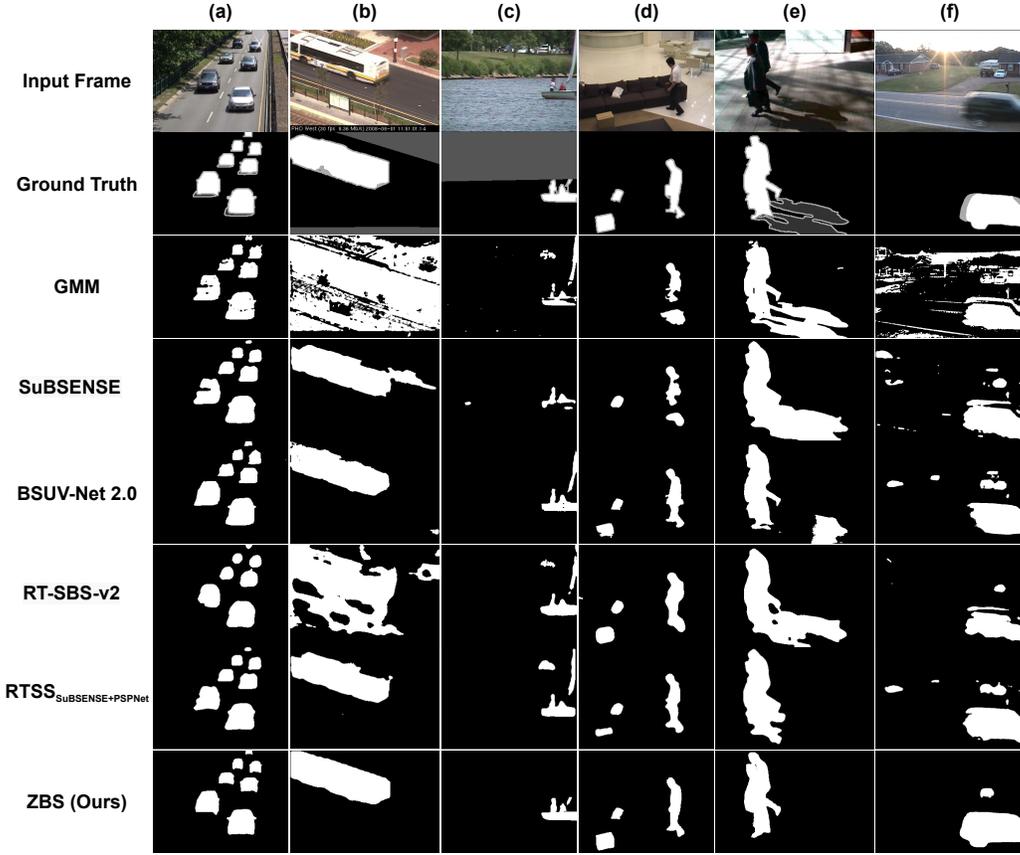


Figure 5. Comparison of the visual results on various scenarios from CDnet 2014 dataset. Except for our method, other segmentation results are quoted in RTSS [43]. From left to right: (a) Scenario "highway" from the *Baseline* category. (b) Scenarios "boulevard" from the *Camera Jitter* category. (c) Scenario "boats" from the *Dynamic Background* category. (d) Scenario "sofa" from the *Intermittent Object Motion* category. (e) Scenario "peopleInShade" from the *Shadow* category. (f) Scenario "continuousPan" from the *PTZ* category.

far outperforms all unsupervised BGS methods and is better than the state-of-the-art supervised method FgSegNet [24]. In the appendix, we also add experiments to demonstrate the good performance in other complex scenes such as night light, camouflaged foreground, *etc.*

$$\text{FPR-S} = \text{nbShadowError} / \text{nbShadow} \quad (5)$$

where nbShadowError is the number of times a pixel is labeled as shadow in Ground Truth but detected as a moving object. nbShadow is the total number of pixels labeled as a shadow in Ground Truth for a video or category.

Table 3. The FPR-S and nbShadowError of different BGS methods.

Method	nbShadowError				FPR-S
	busStation	peopleInShade	bungalows	cubicle	
FgSegNet [24]	2383	12866	5375	580	0.0042
BSUV-Net 2.0 [37]	23149	564989	982943	111438	0.2506
SuBSENSE [34]	315658	854157	1705793	391569	0.5996
SemanticBGS [9]	169426	782489	730668	33137	0.3018
RT-SBS-v2 [13]	28530	457467	566642	52746	0.1717
ZBS (Ours)	964	1892	10403	390	0.0019

5. Conclusion

In this paper, we propose a novel background subtraction framework, ZBS, consisting of three components: all-instance detection, instance-level background modeling, and foreground instance selection. Experiments on the CDnet 2014 dataset show the algorithm's effectiveness. Compared with other BGS methods, our method achieves state-of-the-art performance among all unsupervised BGS methods and even outperforms many supervised deep learning algorithms. ZBS detects most real-world categories without pre-defined foreground categories, producing accurate foreground edges and reducing false detections.

Our future work is leveraging instance-level information more effectively and compressing the all-instance detector for better efficiency.

Acknowledgement. This work was supported by National Key R&D Program of China under Grant No.2021ZD0110403. This work was also supported by National Natural Science Foundation of China under Grants 61976210, 62176254, 62006230, 62002357, and 62206290.

References

- [1] Fateme Bahri and Nilanjan Ray. Dynamic background subtraction by generative neural networks. In *2022 18th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–8. IEEE, 2022. 2
- [2] Mohammed Chafik Bakkay, Hatem A Rashwan, Houssam Salmame, Louahdi Khoudour, D Puig, and Yassine Ruichek. BSCGAN: Deep background subtraction with conditional generative adversarial networks. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 4018–4022. IEEE, 2018. 2
- [3] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 384–400, 2018. 3
- [4] Olivier Barnich and Marc Van Droogenbroeck. ViBe: A universal background subtraction algorithm for video sequences. *IEEE Transactions on Image processing*, 20(6):1709–1724, 2010. 1
- [5] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Uprocft. Simple online and realtime tracking. In *2016 IEEE international conference on image processing (ICIP)*, pages 3464–3468. IEEE, 2016. 4
- [6] Simone Bianco, Gianluigi Ciocca, and Raimondo Schettini. Combination of video change detection algorithms by genetic programming. *IEEE Transactions on Evolutionary Computation*, 21(6):914–928, 2017. 6, 7
- [7] Thierry Bouwmans. Traditional and recent approaches in background modeling for foreground detection: An overview. *Computer science review*, 11:31–66, 2014. 1
- [8] Thierry Bouwmans, Sajid Javed, Maryam Sultana, and Soon Ki Jung. Deep neural network concepts for background subtraction: A systematic review and comparative evaluation. *Neural Networks*, 117:8–66, 2019. 2
- [9] Marc Braham, Sebastien Pierard, and Marc Van Droogenbroeck. Semantic background subtraction. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages 4552–4556. Ieee, 2017. 2, 7, 8
- [10] Marc Braham and Marc Van Droogenbroeck. Deep background subtraction with scene-specific convolutional neural networks. In *2016 international conference on systems, signals and image processing (IWSSIP)*, pages 1–4. IEEE, 2016. 2
- [11] Yingying Chen, Jinqiao Wang, Bingke Zhu, Ming Tang, and Hanqing Lu. Pixelwise deep sequence learning for moving object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(9):2567–2579, 2017. 2
- [12] Pojala Chiranjeevi and S Sengupta. New fuzzy texture features for robust detection of moving objects. *IEEE Signal Processing Letters*, 19(10):603–606, 2012. 1
- [13] Anthony Cioppa, Marc Van Droogenbroeck, and Marc Braham. Real-time semantic background subtraction. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 3214–3218. IEEE, 2020. 2, 6, 7, 8
- [14] Ahmed Elgammal, Ramani Duraiswami, David Harwood, and Larry S Davis. Background and foreground modeling using nonparametric kernel density estimation for visual surveillance. *Proceedings of the IEEE*, 90(7):1151–1163, 2002. 1
- [15] Belmar Garcia-Garcia, Thierry Bouwmans, and Alberto Jorge Rosales Silva. Background subtraction in real applications: Challenges, current models and future directions. *Computer Science Review*, 35:100204, 2020. 1
- [16] Nil Goyette, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, and Prakash Ishwar. Changedetection.net: A new change detection benchmark dataset. In *2012 IEEE computer society conference on computer vision and pattern recognition workshops*, pages 1–8. IEEE, 2012. 6
- [17] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 3, 4
- [18] Enqiang Guo, Xinsha Fu, Jiawei Zhu, Min Deng, Yu Liu, Qing Zhu, and Haifeng Li. Learning to measure change: Fully convolutional siamese metric networks for scene change detection. *arXiv preprint arXiv:1810.09111*, 2018. 2
- [19] Sahin Isik, Kemal Özkan, Serkan Günal, and Ömer Nezh Gerek. SWCD: a sliding window and self-regulated learning-based background updating method for change detection in videos. *Journal of Electronic Imaging*, 27(2):023002, 2018. 6, 7
- [20] S. Jabri, Zoran Duric, Harry Wechsler, and Azriel Rosenfeld. Detection and location of people in video images using adaptive fusion of color and edge information. *International Conference on Pattern Recognition*, 2000. 1
- [21] Kyungnam Kim, Thanarat H Chalidabhongse, David Harwood, and Larry Davis. Background modeling and subtraction by codebook construction. In *ICIP04*, pages 3061–3064, 2004. 1
- [22] Sang-ha Lee, Gyu-cheol Lee, Jisang Yoo, and Soonchul Kwon. Wisenetmd: Motion detection using dynamic background region analysis. *Symmetry*, 11(5):621, 2019. 6, 7
- [23] Zhihui Li, Lina Yao, Xiaoqin Zhang, Xianzhi Wang, Salil S. Kanhere, and Huaxiang Zhang. Zero-shot object detection with textual descriptions. *national conference on artificial intelligence*, 2019. 3
- [24] Long Ang Lim and Hacer Yalim Keles. Foreground segmentation using convolutional neural networks for multiscale feature encoding. *Pattern Recognition Letters*, 112:256–262, 2018. 1, 2, 6, 7, 8
- [25] Long Ang Lim and Hacer Yalim Keles. Learning multi-scale features for foreground segmentation. *Pattern Analysis and Applications*, 23(3):1369–1380, 2020. 7
- [26] Kevin Lin, Shen-Chi Chen, Chu-Song Chen, Daw-Tung Lin, and Yi-Ping Hung. Abandoned object detection via temporal consistency modeling and back-tracing verification for visual surveillance. *IEEE Transactions on Information Forensics and Security*, 2015. 11
- [27] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014. 3

- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 3, 4
- [29] Shafin Rahman, Salman Khan, and Nick Barnes. Improved visual-semantic alignment for zero-shot object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 11932–11939, 2020. 3
- [30] Gani Rahmon, Filiz Bunyak, Guna Seetharaman, and Kannappan Palaniappan. Motion U-Net: Multi-cue encoder-decoder network for motion segmentation. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 8125–8132. IEEE, 2021. 7
- [31] Dimitrios Sakkos, Heng Liu, Jungong Han, and Ling Shao. End-to-end video background subtraction with 3d convolutional neural networks. *Multimedia Tools and Applications*, 77(17):23023–23041, 2018. 2
- [32] Marcos CS Santana, Leandro Aparecido Passos, Thierry P Moreira, Danilo Colombo, Victor Hugo C de Albuquerque, and Joao Paulo Papa. A novel siamese-based approach for scene change detection with applications to obstructed routes in hazardous environments. *IEEE Intelligent Systems*, 35(1):44–53, 2019. 2
- [33] Bruno Sauvalle and Arnaud de La Fortelle. Autoencoder-based background reconstruction and foreground segmentation with background noise estimation. *arXiv preprint arXiv:2112.08001*, 2021. 2
- [34] Pierre-Luc St-Charles, Guillaume-Alexandre Bilodeau, and Robert Bergevin. SuBSENSE: A universal change detection method with local adaptive sensitivity. *IEEE Transactions on Image Processing*, 24(1):359–373, 2014. 1, 6, 7, 8
- [35] Pierre-Luc St-Charles, Guillaume-Alexandre Bilodeau, and Robert Bergevin. A self-adjusting approach to change detection based on background word consensus. In *2015 IEEE winter conference on applications of computer vision*, pages 990–997. IEEE, 2015. 1, 6, 7
- [36] Chris Stauffer and W Eric L Grimson. Adaptive background mixture models for real-time tracking. In *Proceedings. 1999 IEEE computer society conference on computer vision and pattern recognition (Cat. No PR00149)*, volume 2, pages 246–252. IEEE, 1999. 1
- [37] M Ozan Tezcan, Prakash Ishwar, and Janusz Konrad. BSUV-Net 2.0: Spatio-temporal data augmentations for video-agnostic supervised background subtraction. *IEEE Access*, 9:53849–53860, 2021. 2, 6, 7, 8
- [38] Ozan Tezcan, Prakash Ishwar, and Janusz Konrad. BSUV-Net: A fully-convolutional neural network for background subtraction of unseen videos. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2774–2783, 2020. 2
- [39] Han Vanholder. Efficient inference with tensorrt. In *GPU Technology Conference*, volume 1, page 2, 2016. 7
- [40] Yi Wang, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, Yannick Benezeth, and Prakash Ishwar. CDnet 2014: An expanded change detection benchmark dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 387–394, 2014. 5
- [41] Yi Wang, Zhiming Luo, and Pierre-Marc Jodoin. Interactive deep learning method for segmenting moving objects. *Pattern Recognition Letters*, 96:66–75, 2017. 2, 6, 7
- [42] Yizhong Yang, Jiahao Ruan, Yongqiang Zhang, Xin Cheng, Zhang Zhang, and Guangjun Xie. STPNet: A spatial-temporal propagation network for background subtraction. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):2145–2157, 2021. 2, 6, 7
- [43] Dongdong Zeng, Xiang Chen, Ming Zhu, Michael Goesele, and Arjan Kuijper. Background subtraction with real-time semantic segmentation. *IEEE Access*, 7:153869–153884, 2019. 2, 6, 7, 8
- [44] Hongxun Zhang and De Xu. Fusing color and texture features for background model. In *Fuzzy Systems and Knowledge Discovery: Third International Conference, FSKD 2006, Xi'an, China, September 24-28, 2006. Proceedings 3*, pages 887–893. Springer, 2006. 1
- [45] Chenqiu Zhao, Kangkang Hu, and Anup Basu. Universal background subtraction based on arithmetic distribution neural network. *IEEE Transactions on Image Processing*, 31:2934–2949, 2022. 2
- [46] Xu Zhao, Yingying Chen, Ming Tang, and Jinqiao Wang. Joint background reconstruction and foreground segmentation via a two-stage convolutional neural network. In *2017 IEEE International Conference on Multimedia and Expo (ICME)*, pages 343–348. IEEE, 2017. 2
- [47] Wenbo Zheng, Kunfeng Wang, and Fei-Yue Wang. A novel background subtraction algorithm based on parallel vision and bayesian gans. *Neurocomputing*, 394:178–200, 2020. 2, 7
- [48] Dongxiang Zhou and Hong Zhang. Modified gmm background modeling and optical flow for detection of moving objects. In *2005 IEEE international conference on systems, man and cybernetics*, volume 3, pages 2224–2229. IEEE, 2005. 1
- [49] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 350–368. Springer, 2022. 2, 3, 4, 5

A. Instance-level foreground detection

Unlike previous methods, our method builds an instance-level background model. Therefore, ZBS can achieve instance-level foreground detection. Figure 6 shows the difference between binary foreground detection and instance-level foreground detection. Figure 6b shows that ZBS can detect moving foreground of different granularities, including person, backpack, shoe, beanie, etc., and can correctly classify stationary subway and crossbar as background.

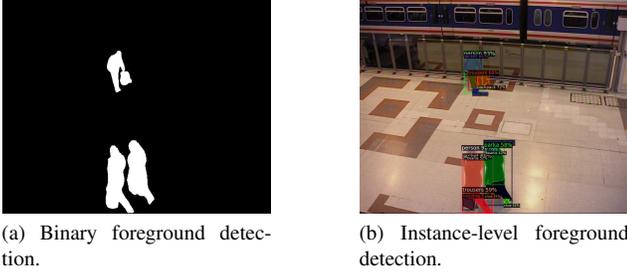


Figure 6. The binary and instance-level foreground detection of ZBS. Our method can detect the moving foreground of different granularities.

B. Abandoned object detection

Abandoned object detection in video surveillance is critical for ensuring public safety and is a crucial component of Intelligent Monitoring. This task presents a challenge, as the categories of abandoned objects are highly diverse and difficult to learn through traditional supervised training methods. Traditional background subtraction techniques often prove insufficient in addressing this issue. Our proposed method, however, offers a solution by incorporating a stronger semantic discernment and instance-level background model, resulting in effective detection of abandoned objects.

To adapt to new tasks, we have added a new rule that considers both motion information and the relationships between instances. If an object exhibits isolated, static behavior or moves independently after previously moving in sync with categories such as a person or car, the instance is deemed to be an abandoned object. This straightforward semantic rule has proven to be effective in diverse environments. We have conducted thorough experiments on the public datasets PETS2006 and ABODA, as well as a non-public traffic abandoned object detection dataset known as TADA.

B.1. PETS2006

The PETS2006 dataset includes sequences from seven different scenes. Each sequence contains an abandonment event except for the third event. We evaluate all seven sequences and our method successfully detects the abandoned objects for the entire PETS2006 dataset without any false alarms. As shown in Figure 7, the results from the PETS2006 dataset demonstrate the efficacy of our ap-

B.2. ABODA

The **AB**andoned **O**bjects **D**ataset (ABODA) [26] contains 11 sequences that present a range of challenging scenarios for abandoned object detection, including crowded scenes, changes in illumination, night-time detection, and both indoor and outdoor environments. Figure 8 displays

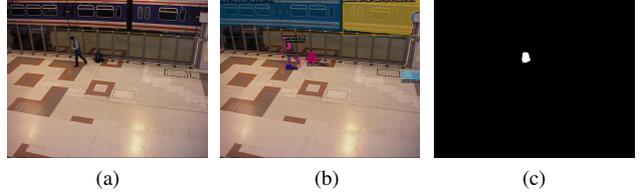


Figure 7. The detection results of the PETS2006. (a) is the original frame. (b) is the all-instance detection results. (c) is the abandoned object detection results of our method.

the results from the *video1.avi* sequence in ABODA.



Figure 8. The detection results of the ABODA. (a) is the original frame in the video. (b) is the all-instance detection results. (c) is the abandoned object detection results of our method.

B.3. TADA

The **T**raffic **A**bandoned object detection **D**ataset (TADA) is a household traffic abandoned object detection dataset that comprises 20 sequences, 14 of which contain traffic abandoned objects. These objects typically consist of various types of traffic litter, such as plastic bags, which have diverse appearances and shapes and are usually carried by the wind. This presents a formidable challenge for abandoned object detection. Figure 9 displays the results from the TADA dataset.

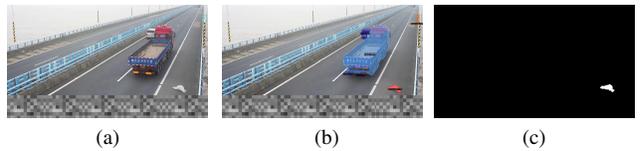


Figure 9. The detection results of the TADA. (a) is the original frame in the video. (b) is the all-instance detection results. (c) is the abandoned object detection results of our method.