*Article*

# A lightweight video detection framework based on information theory and machine learning

**Qi Liu[1], Guangtai Ding[2]**

[1]  School of Computer Engineering and Science, Shanghai University, Shanghai, China; lq1326@i.shu.edu.cn
[2]  School of Computer Engineering and Science, Shanghai University, Shanghai, China; gtding@i.shu.edu.cn
*  Correspondence: gtding@shu.edu.cn

**Abstract:** In recent years, many algorithms based on end-to-end deep networks have been proposed to deal with the target detection problem of videos. However, the deep network models usually consume a lot of computing resources during the procedure of analysis of videos with complex dynamic backgrounds. In this paper, a new method of object detection based on information theory is presented. Firstly, each frame in a video is converted into an effective information map by using the Harris corner detection method. Secondly, the sensitive areas in the frame are extracted by using the context information and the effective information maps of the consecutive video frames. The sensitive areas in the video frame are the candidate areas where the target objects would be appeared at high probabilities. Thirdly, the information entropy features of each sensitive area are extracted to form the feature matrix, based on which, an SVM model is trained for selecting the target areas from the sensitive areas. Finally, the locations of the objects are detected based on the target areas in the video with a complex dynamic background. As a lightweight video detection framework, the method presented in this paper can save a lot of computing resources. Experimental results show that this method can achieve good results in the benchmark of CDnet 2014.
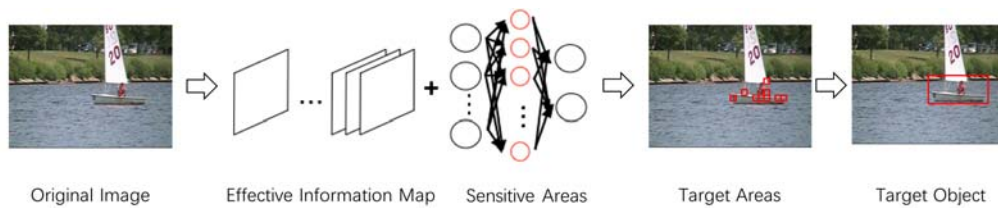
## 1. Introduction

In the last decade, the video processing technology of target detection and tracking has made great progress. Without considering the computational cost of hardware, the most common and effective methods usually depend on the deep neural networks. The end-to-end convolutional models often have good performance in speed and accuracy [1]. Actually, the models are mainly designed for the target detection of still images. In the video processing, these models need to be optimized with the contextual information. When some interfering factors, such as illumination change, motion blur or video jitter, appear in the video, the accuracy of recognition of them would decline sharply.

The speed of video processing rests with the positioning time of suspected target areas in the video frame [2]. The computational cost of video processing depends on the complexity of the deep network model [3]. According to the regularity of the background change and the finiteness of the target category, the video processing model can be designed and modified specially for some certain scenes. For example, the NoScope model is made up of difference detector, specialized model and reference neural network [4]. Due to the special structure and the search method of the model, it can both improve processing speed and reduce operation cost. But after the complex training, the NoScope model is highly dependent on specific scene and target.

The video processing requires avoiding the use of complex deep neural network to reduce hardware cost. At the same time, it also requires exploiting the data redundancy between consecutive frames to reduce software (computation) costs. Zhu [5] uses sparse feature propagation to save these

46  cost. These propagated features are only calculated on the key frames in the video. But this kind of
47  algorithm still needs to be improved because there is not a proper method to dynamically select the
48  key frames from all video frames.
49



Original Image          Effective Information Map          Sensitive Areas          Target Areas          Target Object

50
51                                    **Figure 1.** The flow chart of method.

52          In this paper, a method based on information theory is proposed to detect the video object at
53  low cost. The flow chart of the method is shown in Figure 1. In a video with complex and changeable
54  background, it is very difficult to distinguish the target from the noise only by the motion trajectory.
55  But the change of local information entropy between consecutive frames is effective to identify the
56  object. Based on these features, it is possible to design a relatively simple model for noise reduction.
57  Thus, the hardware which supports this model to run can save a lot of cost. In addition, the sensitive
58  areas in the video frame are more likely to contain foreground targets. So it can avoid a lot of time-
59  consuming operations by doing intensive calculation for sensitive areas rather than that for the whole
60  frame image. The objective of experiments in this paper is to detect objects in the video with dynamic
61  background. And the definition of target detection in this paper is only to mark object positions, not
62  to classify the objects. Comprehensive experiments show that the method based on information
63  theory achieves high accuracy and significant performance.

64  **2. Previous Work**

65          In the past period, there has been a lot of discussion and research on target detection and object
66  classification in the field of video processing. On the basis of the image processing algorithms, a large
67  number of video processing algorithms has been proposed and improved. In practice, many
68  experiments show that combining traditional feature extraction methods (like SIFT [6] and HOG [7])
69  and machine learning methods can solve the problem of object detection and classification well in the
70  field of image processing [8,9]. In DPM model, the object is divided into multiple parts for extracting
71  HOG features separately [10]. These HOG features are proved to be useful and efficient in pedestrian
72  detection. Due to the lower computational cost and higher detection accuracy, deep neural network
73  becomes an important method of computer vision processing in recent years. Current evidence
74  suggests that the deep learning algorithm has an ability of approximate human beings in the field of
75  image classification and object detection [11,12].
76          When all of the frames in the video are processed in the same way, the applications of deep
77  network in image processing can be directly converted to algorithms for solving video processing
78  tasks. Without considering the cost and the speed, the task of object recognition and segmentation in
79  video processing can be solved well by improving the deep network algorithm (like CNN [13]) in
80  image processing [14,15]. But as the requirements for speed and accuracy are gradually improved,
81  more models are specially designed as end-to-end deep network models [16,17]. These models can
82  avoid many problems about parameter optimization during the training process. Meanwhile they
83  can also speed up by directly outputting the category and the location of object.
84          The end-to-end network model is considered as the key point of improving the speed of
85  recognition and positioning. In the object recognition task of computer vision, the YOLO model
86  abandons the thought of regional pre-processing. And it truly realizes the application of end-to-end
87  network model [18,19]. Now, without considering the cost of hardware, the target object in video
88  processing task can be identified and tracked in real time, such as SSD model [20]. Zhu [5] proposes
89  the idea of transferring the convolutional feature map of sparse key frame to other video frames. This
90  idea can take advantage of the context information well in the video stream and accelerate the

91   operation process of the whole model by flow calculation. This method is specially designed for
92   solving video processing tasks.
93       The use of end-to-end deep network needs to consume a lot of costs, such as manual marking
94   cost and hardware cost. So before the end-to-end deep learning network becomes popular, R-CNN
95   (regional convolutional neural network) is the primary method to solve the object recognition
96   problem in the field of video processing. Compared to a complex deep network, R-CNN may be run
97   with lower hardware cost but will be followed by the inefficiency of computation. The method
98   proposed in this paper tries to solve this contradiction. It improves the speed of video processing by
99   using sensitive area screening. At the same time, it saves hardware cost by using simple machine
100  learning model.
101      The function of information entropy is to quantify the information change process of the local
102  area. The information entropy of one image is actually the expected value of all information saved in
103  this image. Therefore, the change of information entropy in continuous frames is one special form of
104  information gain. The information gain is often used in decision tree algorithm to select
105  characteristics [25]. Actually, the random forest algorithm proposed based on decision tree is still
106  popular in many fields now [26]. These works prove that the information theory is useful for finding
107  the effective characteristics of different classes. Therefore, the ability of information theory in the
108  classification model is not to be questioned. In this paper, the feature matrix used in the training of
109  model is obtained by calculating the information entropy of sensitive area. These feature matrixes
110  can better express the deep characteristics of image and improve the generalization ability of the
111  model.

## 3. Method

113      The method proposed in this paper is to achieve the object detection task in a complex
114  environment. The video processing in complex environment is usually more difficult. Thus a new
115  method in this paper is proposed on the basis of information theory. In the method, the information
116  entropy is utilized to quantify the information change process of the local area. The quantify result
117  can well indicate the different nature between the object and the noise. As shown in Figure 1, the
118  focus of the method is the extraction of sensitive areas in the video frame and then the feature matrix
119  is extracted by information entropy calculation to classify sensitive areas.
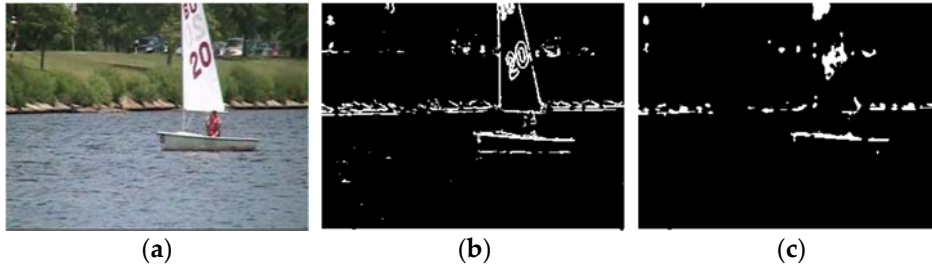
### 3.1. Effective Information Map

121      The effective information represents the valuable information that is used for identifying
122  foreground objects in the frame of video. There is a greater probability of having targets in the area
123  where effective information is assembled in the image. The effective information map can be used to
124  extract sensitive areas of the image, such that the target can be found faster by less computation.
125  Experiments show that about 1000-2000 candidate areas can be found out from a video frame in the
126  pre-process of region-based convolutional networks [14]. However, extracting sensitive areas
127  through effective information map can greatly reduce the number of these candidate areas.
128      In selective search model, the candidate area set is composed of local regions that are merged
129  after image segmentation [2]. Therefore, the angular points and edges that are important during the
130  extracting processing of candidate areas often represent effective information. Since the Harris corner
131  detection algorithm [21] can distinguish the flat regions, the edge regions and the corner regions in
132  the image, this algorithm can be used as the extraction method of effective information.
133      The Harris corner detection algorithm can obtain the response value of the corner information
134  through the transformation of corner response function. The response value of the edge region is
135  negative. The response value of the corner region is a larger positive number and the response value
136  of the relatively flat region is a smaller positive number. Assume that the non-flat region in the video
137  frame is the area containing effective information, so the transformation formula of the effective
138  information map is

139
$$I_{\text{info}} = f(x, y) = \begin{cases} 255, & |dst| > (\alpha \cdot \text{Max}(|dst|)) \\ 0, & |dst| \le (\alpha \cdot \text{Max}(|dst|)) \end{cases}$$

140  $dst$ represents the corner detection result of the video frame. $\alpha$ times of the maximum value of the
141  corner detection result is the threshold to divide the flat region and the non-flat region, i.e., the
142  ineffective information region and the effective information region. Finally, the effective information
143  map can be obtained by the Gaussian filtering and simple morphological processing. As shown in
144  Figure 2, using the Harris corner detection algorithm can filter out a large number of flat regions,
145  such as lake surface, grass lawn and so on. Then, the sensitive area in video frame can be clearly
146  found out from the effective information map obtained through subsequent processing.

147



148
149                    (**a**)                              (**b**)                              (**c**)

150               **Figure 2.** (**a**) Original Image (**b**) Harris Image (**c**) Effective Information Map

151  *3.2. Sensitive Area Extraction*

152       Using the context information of video is very important to improve video processing efficiency.
153  The frame difference and the background difference are common methods of using context
154  information. In these methods, the consecutive frame of video rather than a single frame is processed
155  and calculated as a basic unit. The dynamic background often contains complex information. So the
156  video with dynamic background is generally difficult to be processed by using the background
157  difference method. The frame difference method is characterized by low complexity, fast running
158  speed and strong adaptive ability of dynamic environment. Some noise in the dynamic background
159  can sometimes be mistaken for the foreground object during the processing of the frame difference
160  method. So this paper tries to improve the frame difference method to avoid the influence of noise as
161  far as possible.
162       According to the effective information map rather than the original video frame as the
163  processing unit, the frame difference method can get better result. The method used in this paper is
164  to calculate on the basis of effective information maps of three consecutive frames. The calculation
165  formula is

166          $$D_n(x, y) = [f_n(x, y) - f_{n+1}(x, y) \wedge f_n(x, y)] \vee [f_n(x, y) - f_n(x, y) \wedge f_{n-1}(x, y)]$$

167  $f_n(x, y)$ represents the effective information map of the n'th video frame. The $\vee$ operation in the
168  formula is to calculate the mean value of corresponding pixels in two effective information maps. The
169  $\wedge$ operation in the formula is to reserve the corresponding pixel value of the effective information
170  map of n'th video frame. This method can retain the outline of the moving object well after removing
171  the background area. The threshold processing is required for $D_n(x, y)$ after the difference operation.
172  The threshold value $T_{Otsu}$ is obtained by the Otsu [24] method automatically. Then the influence of
173  light fluctuation is added on this basis to get the optimal threshold value $T_{optimal}$. The optimal
174  threshold calculation formula is

175          $$L_{\text{Diff}} = \frac{1}{N_A} \sum_{(x,y) \in A} (|f_{n+1}(x, y) - f_n(x, y)| + |f_n(x, y) - f_{n-1}(x, y)|)/2$$

176                              $$T_{\text{optimal}} = T_{\text{Otsu}} + \lambda \cdot L_{\text{Diff}}$$

177  $\lambda$ represents the influence factor of light fluctuation in the current environment. Through the
178  threshold processing, the difference image can be obtained finally.
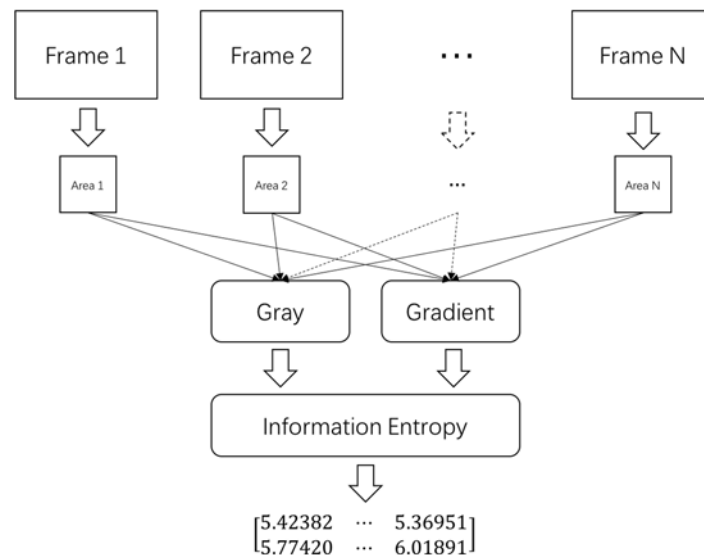
179      The sensitive area is the area where the foreground object is easy to appear in video frame.
180  Sensitive areas can be determined easily in the difference image of consecutive video frames. The
181  difference image obtained after threshold processing is divided into many small regions. These
182  regions have the same size and they are non-overlapping. The number of effective information pixel
183  in each region represents the occurrence possibility of the foreground object. This occurrence
184  possibility is the basis for judging sensitive areas. This paper assumes that the region is judged as
185  sensitive area as long as its occurrence possibility is greater than zero. The acquisition of sensitive
186  areas can greatly save the computing cost of subsequent operations.

187  *3.3. Sensitive Area Screening*

188      The extraction of sensitive area can effectively save the calculation cost in the subsequent
189  operation. The another key of the method in this paper is how to select real target areas from all of
190  the sensitive areas. In the most recent period, a large number of end-to-end machine learning models
191  are used to recognize targets [11]. But not all classification problems need to be solved by the deep
192  network [22]. In the NoScope model, just a simple convolutional neural network is used to recognize
193  objects that belong to a small number of categories [4]. The training of deep network often requires a
194  large number of sample labels. During the running process, the deep network also consumes a lot of
195  hardware resources. In the video which contains a dynamic background, these shortcomings of the
196  deep network are easy to be magnified because of many noise of the background.
197      In the video with complex background, the location recognition of the foreground object is
198  important than the classification of the foreground object. During the consecutive video frames, the
199  information changes of the target area are often progressive. Intuitively, the information changes of
200  the noise area tend to have the characteristics of step change. Therefore, this paper proposes a method
201  that uses information entropy to quantify the process of information change of the area. Then, a
202  simple binary classification model can be trained through these information entropy matrixes. Since
203  the model is based on information change of local area, rather than image area itself, to distinguish
204  the foreground object from the background noise, the model has the characteristics of simple, fast
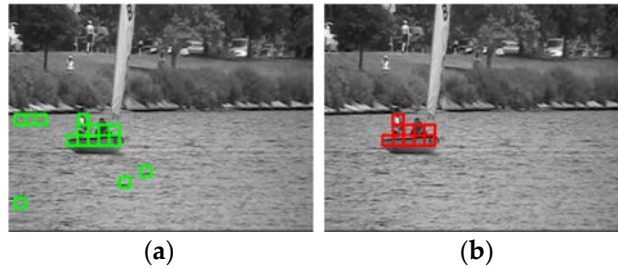205  and strong robust.
206



208      **Figure 3.** The extraction flow chart of feature matrix of information entropy.

209      The extraction process of information entropy feature of the sensitive area is shown in the Figure
210  3. An N dimensional feature vector can be obtained by calculating the information entropy of local
211  areas of the same location in consecutive N frames. But the local areas of video frame are firstly
212  preprocessed by M kinds of algorithms before the calculation of feature vector. Thus, one area in the
213  video frame can be finally represented by a M×N dimensional feature matrix. There are many options
214  for image processing algorithm in preprocessing operation. Two preprocessing algorithms used in

215    this paper are image gray algorithm and image gradient algorithm. No matter which kind of methods,
216    it can be of more generality after the calculation of information entropy.
217



(**a**)                                     (**b**)

**Figure 4.** (**a**) Sensitive Areas (**b**) Target Areas

221    During the training of the machine learning model, one feature matrix is used as the sample data
222    of one sensitive area. These samples are used for training a supervised binary classification model.
223    As shown in Figure 4, this model is used for screening sensitive areas in the video frame, in other
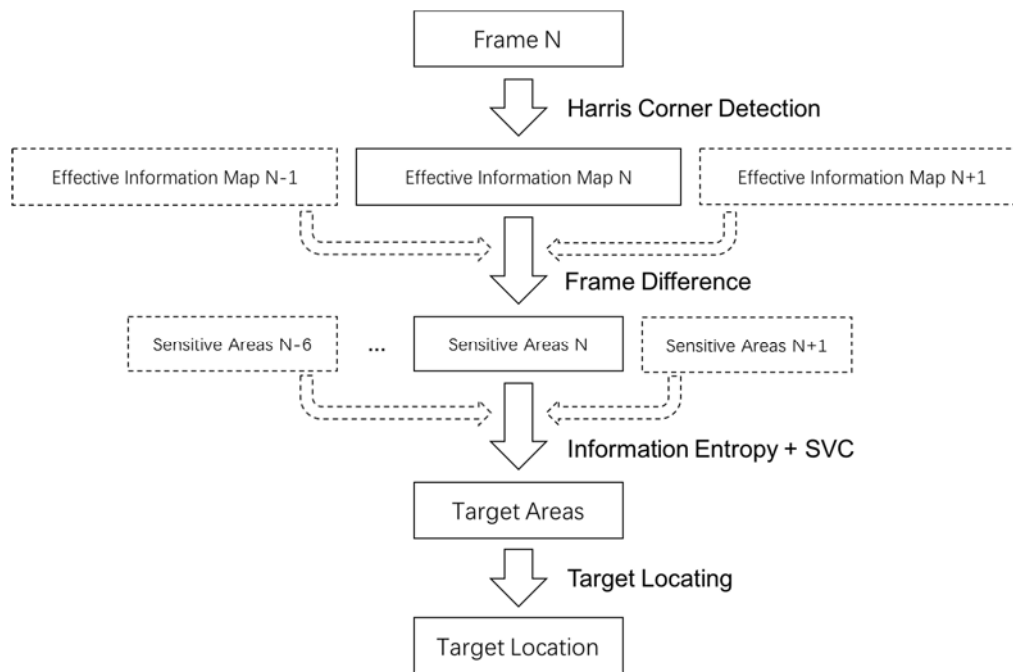224    words, filtering out the noise generated by dynamic background.

225    *3.4. Target Object Locating*

226    The target object is calibrated on the basis of target areas in the video frame. Since the previous
227    operation has divided the whole image into multiple regions, the result image which contains all
228    target areas is resized for the convenience of subsequent processing. The adjustment rule of image
229    size is that each target area corresponds to one pixel in the adjusted image. Through the
230    morphological analysis of the adjusted image, target areas which belongs to the same object are
231    merged into a large target block. After the image is adjusted back to its previous size, the four
232    boundaries of this irregular target block can be utilize to form a complete rectangle. Finally, these
233    rectangular boxes can be used as the calibration results of the target objects in the video frame. In the
234    next moment, they can be used for classification or other more significant work. The method that uses
235    target areas for locating the target object is not complicated, but it can achieve a nice result.

236    **4. Experiment**

237    The experiment data of this paper is from the video dataset in CDnet 2014 [23]. The main test set
238    used in this paper is the video with dynamic background. The ground truth image in dataset only
239    distinguish the foreground object from the background. Thus, the detection task of this paper mainly
240    focuses on the calibration of the target locations, not on the recognition of the target category. One of
241    the main contributions of this paper is to use a more efficient algorithm to achieve the selection work
242    of candidate areas. The implementation of the regional convolutional network usually leads to
243    excessive running cost because of the time-consuming selection work in it. The end-to-end deep
244    network does not need to screen out candidate areas in advance, so it becomes more popular now.
245    The work in this paper has positive significance to use a low-cost simple model for achieving the
246    same effect with an end-to-end deep network.
247    The experimental flow chart of the method proposed in this paper is shown in Figure 5. During
248    the transform process of effective information map, the threshold of Harris corner detection
249    algorithm is 0.01. If the shot environment of video is not very complex, extreme or uncommon, then
250    the influence of this parameter is little. During the extraction and screening of sensitive areas,
251    experiments prove that the optimal pixel size of sensitive area is 11×11. After the sensitive area is
252    determined, the same location area of 8 consecutive frames are selected as data sample source for
253    subsequent classification work. Then, each sensitive area can be transformed into a 2×8 dimensional
254    feature matrix.

**Figure 5.** The experimental flow chart of method.

**Table 1.** The model parameters table.

| Kernel | Degree | Gamma | Error Precision |
|--------|--------|-------|-----------------|
| polynomial | 3 | 1/16 | 0.001 |

**Table 2.** The experimental result table.

| Video Category | Dynamic Background | | Camera Jitter |
|----------------|--------|--------|---------------|
| Video Content | Boats | Highway | Traffic |
| Accuracy | 85.29% | 85.61% | 88.05% |
| Speed(fps) | 64 | 52 | 40 |

In the next period of the experiment, different classifiers are trained for several video environments respectively. However, the training parameter of different classifiers are the same. The supervised binary classification model used in the experiments is SVC model. The parameters of model are shown in Table 1. Three representative series of video are selected as primary experiment data: one contains the dynamic background of water fluctuation, one contains the dynamic background of leaves shaking and the other one contains heavy jitter of the camera. As shown in Table 2, the experimental results prove that the object detection framework based on information entropy can achieve the same ideal effect in different scenes. After running the different classification models with cross-validation, the ROC curves are plotted as shown in Figure 6. It can be proved that the object detection framework proposed in this paper has the strong adaptive ability in different scenes.
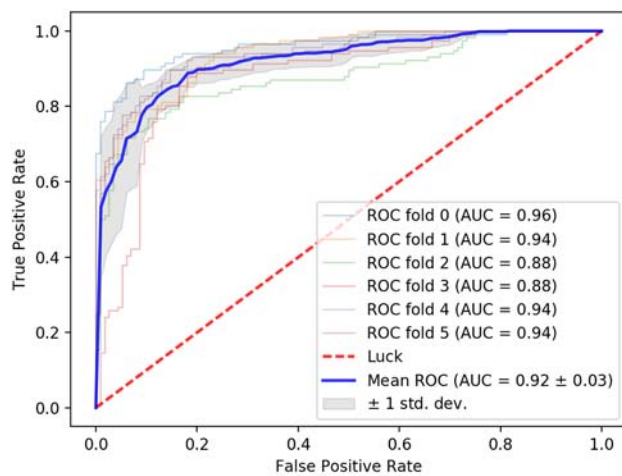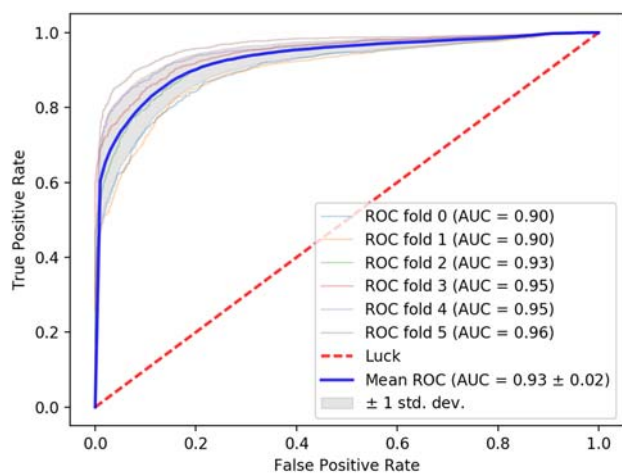
The different classification models based on information theory are trained in different scenes. Then, these models are tested in the corresponding scene. The test results are shown in Figure 7. In the dynamic scene which content is a boat, about 1000 frames in the first half of the video are used in the training of the detection model and about 2000 frames in the second half of the video are used in the test. In Figure 7, the first line is the processing result obtained by using the validation set of video and the second line is the processing result obtained by using the test set of video. If the classification model is trained based on the image itself rather than the information entropy result, the model may get high accuracy in validation set but will be followed by poor generalization ability in the test set.

279    The above experiment proves that the video detection framework based on information theory has
280    strong generalization ability in the same scene.

281

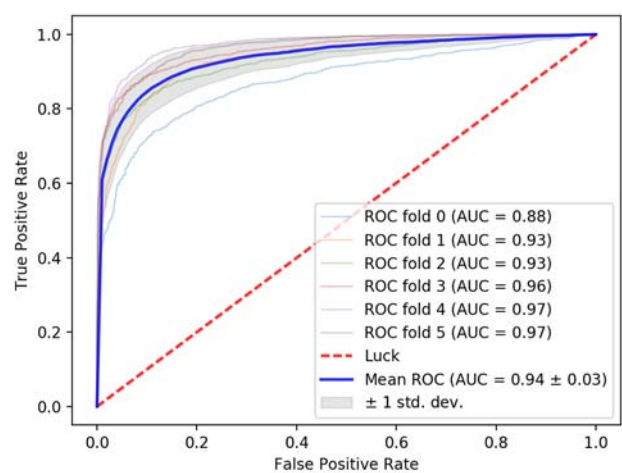282                                                                 (a)

283

284                                                                 (b)
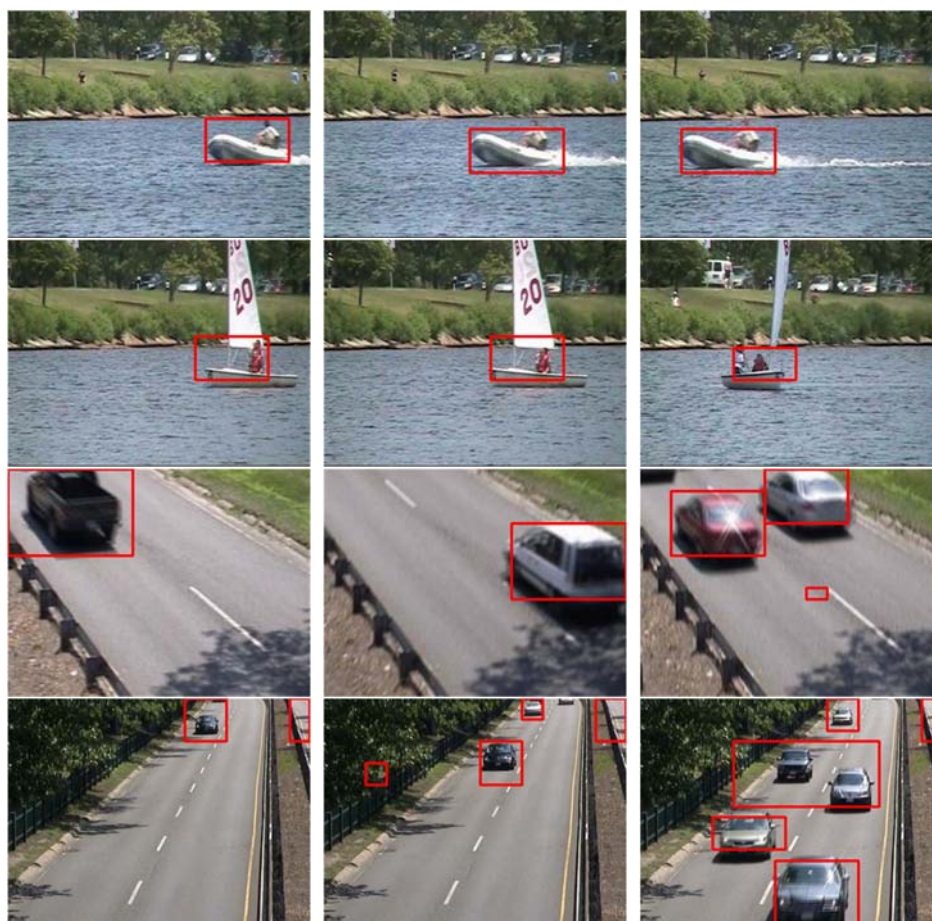
285

286                                                                 (c)

287                    **Figure 6.** The ROC curves. (**a**) Boats (**b**) Highway (**c**) Traffic

**Figure 7.** The test results.

## 5. Conclusion

A new target detection method based on information theory is presented in this paper. Although this method does not have the ability to classify target categories, the method can be treated as a framework for target location detection to save the cost as much as possible. A large number of candidate areas are avoided to be considered by using the context information of the video. At the same time, much hardware cost can be saved by using a lightweight machine learning model. The method of this paper provides a new way, which is different from the end-to-end deep network, to solve the video processing task at low cost.

Although the algorithm has strong generalization ability in the same scene, different models need to be trained respectively for video with different dynamic backgrounds, because the noise has different characteristics in different dynamic background. Thus the future work is to find the common characteristics among the different scenes by using information theory.

## 6. Acknowledgment

## References

1. Huang J, Rathod V, Sun C, et al. Speed/accuracy trade-offs for modern convolutional object detectors[J]. 2016.
2. Uijlings J R, Sande K E, Gevers T, et al. Selective Search for Object Recognition[J]. International Journal of Computer Vision, 2013, 104(2):154-171.
3. Zhu X, Dai J, Yuan L, et al. Towards High Performance Video Object Detection[J]. 2017.

314    4.   Kang D, Emmons J, Abuzaid F, et al. NoScope: Optimizing Neural Network Queries over Video at Scale[J].
315         Proceedings of the Vldb Endowment, 2017, 10(11):1586-1597.
316    5.   Zhu X, Xiong Y, Dai J, et al. Deep Feature Flow for Video Recognition[C]// IEEE Conference on Computer
317         Vision and Pattern Recognition. IEEE, 2017:4141-4150.
318    6.   Lowe D G. Object Recognition from Local Scale-Invariant Features[C]// iccv. IEEE Computer Society,
319         1999:1150.
320    7.   Dalal N, Triggs B. Histograms of Oriented Gradients for Human Detection[C]// Computer Vision and
321         Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on. IEEE, 2005:886-893.
322    8.   Zhu Q, Yeh M C, Cheng K T, et al. Fast human detection using a cascade of histograms of oriented
323         gradients[C]// Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on.
324         IEEE, 2006:1491-1498.
325    9.   Verschae R, Ruiz-Del-Solar J, Correa M. A unified learning framework for object detection and
326         classification using nested cascades of boosted classifiers[J]. Machine Vision & Applications, 2008, 19(2):85-
327         103.
328    10.  Felzenszwalb P F, Girshick R B, Mcallester D, et al. Object detection with discriminatively trained part-
329         based models[J]. Computer, 2014, 47(2):6-7.
330    11.  Russakovsky O, Deng J, Su H, et al. ImageNet Large Scale Visual Recognition Challenge[J]. International
331         Journal of Computer Vision, 2014, 115(3):211-252.
332    12.  Donahue J, Hendricks L A, Guadarrama S, et al. Long-term recurrent convolutional networks for visual
333         recognition and description[M]// AB initto calculation of the structures and properties of molecules /.
334         Elsevier, 2015:85-91.
335    13.  Krizhevsky A, Sutskever I, Hinton G E. ImageNet classification with deep convolutional neural
336         networks[C]// International Conference on Neural Information Processing Systems. Curran Associates Inc.
337         2012:1097-1105.
338    14.  Girshick R, Donahue J, Darrell T, et al. Region-Based Convolutional Networks for Accurate Object
339         Detection and Segmentation[J]. IEEE Transactions on Pattern Analysis & Machine Intelligence, 2016,
340         38(1):142.
341    15.  He K, Zhang X, Ren S, et al. Spatial Pyramid Pooling in Deep Convolutional Networks for Visual
342         Recognition.[J]. IEEE Transactions on Pattern Analysis & Machine Intelligence, 2015, 37(9):1904-16.
343    16.  Girshick R. Fast R-CNN[J]. Computer Science, 2015.
344    17.  Ren S, He K, Girshick R, et al. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal
345         Networks.[J]. IEEE Transactions on Pattern Analysis & Machine Intelligence, 2015, 39(6):1137-1149.
346    18.  Redmon J, Divvala S, Girshick R, et al. You Only Look Once: Unified, Real-Time Object Detection[J].
347         2015:779-788.
348    19.  Redmon J, Farhadi A. YOLO9000: Better, Faster, Stronger[J]. 2016.
349    20.  Liu W, Anguelov D, Erhan D, et al. SSD: Single Shot MultiBox Detector[J]. 2015:21-37.
350    21.  Harris C. A combined corner and edge detector[J]. Proc Alvey Vision Conf, 1988, 1988(3):147-151.
351    22.  Ba L J, Caruana R. Do Deep Nets Really Need to be Deep?[J]. Advances in Neural Information Processing
352         Systems, 2013:2654-2662.
353    23.  Y. Wang, P.-M. Jodoin, F. Porikli, J. Konrad, Y. Benezeth, and P. Ishwar, CDnet 2014: An Expanded Change
354         Detection Benchmark Dataset, in Proc. IEEE Workshop on Change Detection (CDW-2014) at CVPR-2014,
355         pp. 387-394. 2014.
356    24.  Otsu N. A Threshold Selection Method from Gray-Level Histograms[J]. IEEE Transactions on Systems Man
357         & Cybernetics, 2007, 9(1):62-66.
358    25.  Quinlan J R. C4.5: programs for machine learning[J]. 1993, 1.
359    26.  Breiman L. Random Forests[J]. Machine Learning, 2001, 45(1):5-32.