

Application for Detecting Small Objects Capable of Recognizing Text from A Diverse Range of Image Perspectives

Versha Sharma, Rahul Gupta

Department of Computer Science and engineering,
S.R Institute of management & Technology, Lucknow, India,
sversha0987@gmail.com

Abstract: In previous years, advancements in smart city technology have significantly enhanced object detection systems through the adoption of various innovative deep learning models. Deep learning has surpassed traditional computer vision techniques in performance. Recently, automated segmentation models have integrated feature extraction concepts, suggesting diverse feature extractions on images. These models typically employ texture classification and color detection to forecast neighboring pixel labels, with classification validating the features. Parameter tuning in these models typically relies on signal processing strategies. Many researchers have leveraged optimized parameter sets tailored to specific datasets, leading to a significant increase in accuracy. However, these models often lack scalability across different datasets. Introducing a novel hybrid technique combining wavelet and watershed transform for small object detection and text extraction, our approach overcomes this limitation. Our hybrid model demonstrates scalability across various image types, exhibiting improved accuracy even in the presence of diverse noise levels when applied to different photos.

Keywords: Object detection, Masking DWT, watershed, small object, Morphological.

1. Introduction:

An object identification system finds real-world pairings using photographs of the world using an a priori defined object model [1, 2]. The demand for object recognition systems is provisos to ever-growing volume of digital photographs in all public as well as private databases. Humans can do object identification extremely simply and readily, while robots find it quite challenging. Object recognition technologies are critical for robots to achieve increased levels of autonomy [3]. It is a popular field of robotics study that uses computer vision (CV) with machine learning (ML). Drones are increasingly getting employed as robotic devices. The goal of this study is how present object identification algorithms as well as methods may be applied to drone image information. When comparing to certain other robots [4, such as a wheeled robot], another of the benefits of utilizing a drone to identify items in a picture will be that the drone can travel closer to objects. Nevertheless, uses a top - down perspective angels [5] as well as the difficulty of combining using a deep learning system for compute-intensive processes [6] pose challenges regarding UAVs.

Whenever a drone navigates an area in discovery of items, it is advantageous for drone to have as much visibility as feasible [7, 8]. Images captured by UAVs or drones, on the other hand, varied significantly from those captured by a conventional camera. As a result, object identification techniques employed on "normal" photos cannot be presumed to work effectively on photographs captured by drones. Previous research has found that photos obtained by a drone are frequently different from those accessible for training that is frequently shot with a hand-held camera. Study helps to identify in drone information could be tough owing to camera's placement [9, 10] comparing to images captured by a human, based on sort of pictures the network is trained on. Object recognition is an automated image identification procedure dependent on statistical as well as geometric aspects that also demands precise categorization of items in pictures and moreover involves precise estimation of objects.



Figure 1: Small object dataset.

The efficiency of object categorization as object positioning are critical measures of model identification performance. Object recognition is utilized in a variety of applications, including intelligent surveillance, military object identification, UAV navigation, unmanned vehicle navigation, as well as intelligent transportation. The existing model, nevertheless, fails to recognize objects due to the diversity of the identified items. The complexity of object recognition is made more complex by changing light as well as a complicated background, particularly for objects in a

diverse context. The standard multiscale pyramid approach of classification tasks as well as localization requires extracting the picture's statistical multistate data as well as subsequently classifying image with classifier [1–3]. It is tough any but rather many characteristics describe things that don't create strong classification model since various types of photos are defined by distinct characteristics. Those models were unable to identify the items, particularly when there were many recognized objects in a single picture.

2. Related Work:

In the realm of machine intelligence, object recognition is constantly a hot issue. The traditional sliding window recognition approach necessitates the decomposition of pictures in multi-scale pictures. Typically, a single picture is divided into a slew of sub-windows with tens of millions of distinct positions as well as sizes. A classifier is then used by the system to identify whether observed entity is present in every window. Procedure is wasteful since this necessitates a thorough search. Furthermore, various classifiers have an impact on object identification efficiency. Classifiers are developed as per various types of observed in sequence for them to generate a strong classifier. For instance, the Harr feature paired using Adaboosting classifier [14] makes face identification possible. We employ HOG features (Histogram of Gradients) paired using support vector machine [15] for pedestrian recognition, as well as HOG feature mixed through DPM (Deformable Part Model) [16, 17] for common object recognition. Furthermore, if a picture has a large number of various types of recognized items, such entities will be unnoticed by computers..

Since 2014, when Hinton utilize deep learning towards win ImageNet competition with the highest success rate, deep learning becomes a popular method for detecting objects. The object recognition model dependent on deep learning is separated into two classes: the first one is focused on region suggestions [18–20], like RCNN [4], SPP-Net [7], Fast-RCNN [5], Faster-RCNN [6], and R-FCN [8], as well as is the most extensively utilized. The other technique, like YOLO [21] as well as SSD [22], does not employ area suggestions as well as instead identifies the objects immediately.

For the first technique, the model selects RoIs first before detecting them; that is, several RoIs are created using chosen searches [23], edge box [24], or RPN [25]. The system subsequently uses CNN to extract features for every RoI, identifies items using classifiers, as well as determines the position of discovered objects. RCNN [4] generates around 200 RoIs for every image using selective search [23] as well as then isolates as well as classifies convolution characteristics of 200 RoIs, correspondingly. Since there are a lot of people in such RoIs overlapping portions, identification is inefficient due to the enormous number of consecutive computations. For such a situation, SSP-net [7] as well as Fast-RCNN [5] provides mutual RoIs function. The approaches recover just a CNN characteristic from the entire source picture, while so every RoI's characteristic is retrieved separately as from CNN feature using the RoI pooling process. As a result, the level of work required to retrieve every RoI's feature is distributed. This approach lowers the

2000-times CNN operations in RCNN to single CNN operations, considerably improving computational efficiency.

3. Methodology:

With a proportional implementation of the human impression of the intrigue district, the study aims to separate the photographs into distinct groups. These images may be uncontrolled organ division images that have been filtered to meet the ROI presumption criterion. We'll use a two-arrangement strategy. The strategy lobe images division will be implemented, enabling the processing of textured and non-textured photos simultaneously. Sorting the double tree wavelet change coefficients is the initial step towards identifying textured highlights in the image. Subsequently, middle sifting will be linked to reduce surface district ambiguities at the margins of image objects.

Watershed change will be connected to obtain the underlying division once the computed surface component has been used to calculate the space created slope capacity.

Utilizing the cost based on a weighted average effort for district parcelling, the second stage sectioned lands gained by the watershed change are connected to significant locations with comparable qualities.

Image DWT2:

Single-level discrete 2-D wavelet transform

Syntax-

$[cA, cH, cV, cD] = \text{dwt2}(X, 'wname')$

$[cA, cH, cV, cD] = \text{dwt2}(X, Lo_D, Hi_D)$

$[cA, cH, cV, cD] = \text{dwt2}(\dots, 'mode', MODE)$

Description

The `dwt2` function does two-dimensional wavelet decomposition over single level. Comparing the technique with `wavedec2`, that could be better appropriate for your needs. Decomposition is performed either using a specific wavelet ('wname'; see `wfilters` for additional details) or specific wavelet decomposition filters (Lo D as well as Hi D) that you provide.

$[cA, cH, cV, cD] = \text{dwt2}(X, 'wname')$ Approximations coefficients matrix `cA`, as well as details coefficients matrices `cH`, `cV`, as well as `cD` (horizontal, vertical, as well as diagonal, correspondingly), are computed using wavelet decomposition of input matrix `X`. Wavelet identity is included in 'wname' character vector.

$[cA, cH, cV, cD] = \text{dwt2}(X, Lo_D, Hi_D)$ Using wavelet decomposition filters you provide, calculates two-dimensional wavelet decomposition as shown above.

- `Lo_D` is decomposition low-pass filter.
- `Hi_D` is decomposition high-pass filter.

`Lo_D` as well as `Hi_D` should be of equal size.

Let `sx = size(X)` as well as `lf = length of filters`;

then $\text{size}(cA) = \text{size}(cH) = \text{size}(cV) = \text{size}(cD) = sa$ where $sa = \text{ceil}(sx/2)$,

If DWT extensions mode is periodization, for additional forms of extension,

$sa = \text{floor}((sx+lf-1)/2)$.

$[cA, cH, cV, cD] = \text{dwt2}('mode', MODE)$ Wavelet decomposition is computed using `MODE` extensions mode which you select.

MODE is character vector that contains extension mode you want to use.

A good example of appropriate application is
`[cA,cH,cV,cD] = dwt2(x,'dbl1','mode','sym');`

Algorithm:

1. Image obtaining:Analyze biological picture (I) as well as resize it before selecting it locale of intrigue that will be edited. The digital image is in the binary format saved in the binary format. The image is read as the 8 bit integer data and saved in the matrix format. In this code the images e saved by different names and on inserting a specific name the image is imported and saved in matrix data.

2. 1 level Image DWT: Accomplish principal level 2d DWT on the image and acquire the estimation segment (An) of the changed information.

$$[A \text{ DH DV DD}] = \text{DWT}(I)$$

The I is the image matrix and passed through 2 dimensional filters acting along x and y direction along the rows and columns. The DWT applies the filtering to create the high pass and low pass frequency components of image. The low pass component are taken as approximated components saved as A and the high frequency components are DH known as straight details ,DV are vertical details and DD are the diagonal details.

3. 2 level Image DWT: Accomplish second level 2d DWT on image as well as acquire following estimation segment (A1) of changed information.

$$[A1 \text{ DH1 DV1 DD1}] = \text{DWT}(A)$$

The first level DWT gives Components with high as well as low frequencies. Low frequency mechanisms are large in amount as per energy density. They are further broken into low and high frequency. To accomplish this task the low frequency part known as approximation component (A) are further passed through the DWT again and the resultant matrix are called as A1 (aproxx.),DH1,DV1 and DD1 are details part.

4. Approximated Image Reconstruction by 2level IDWT:

The resultant approximated part obtained by 2nd level DWT is in double precision format and called as wavelet coefficients. These approximation coefficients are further added with null value of details component to reconstruct the approximated image in the 8 bit integer format. It is taken as the performing the converse DWT as well as reproduce image by considering just approx segment as well as stifling the DH, DV and DD point by point segment. The details are removed by suppressing them to zero. In this way, the noisy edges and fine lines are removed the image become clear and the tissue morphology become homogenous.

5. Entropy separating: The entropy values represent the information content at different location of image as per the density of pixel intensity variation. The image entropy based separation of high and low entropy regions are segregated to know the information content. Apply entropy based separation to remake approximated image. These entropy values are taken as features to segment the image in different regions.

$$H = - \sum_{i=0}^{255} p_i \log_2 p_i$$

H=entropy of image

p_i =probability of occurrence of a symbol i.

6. Remove little protests: The image consist of dots or circular shapes as openings. The small openings are removed out to minimize ambiguity during the segmentation process. The removal of the undesirable openings having size lower than 10 pixels is performed in this code for making higher accuracy in segmentation.

7. Morphological Processing: The morphological process are performed to bring changes in the shape of image objects. This is performed by two different schemes known as erosion and dialation. The erosion removes the boundary pixels of image objects. It helps in separating two objects, which are closer. The dilation process adds on pixel at the boundary of object. It helps in aggregating object having very thin border lines.

E be a Euclidean space or an integer grid, A a binary image in E , and B a structuring component regarded as a subset of R^d .

The dilation of A by B is defined by

$$A \oplus B = \bigcup_{b \in B} A_b$$

where A_b is the translation of A by b .

Dilation is commutative, also given by

$$A \oplus B = B \oplus A = \bigcup_{a \in A} B_a$$

8. Texture Masking: Textures are the patterns of similar type associated with variation of color and shades and geometrical shape. The biomedical organ consists of tissues of specific color, shade and pattern and they may be easily separated out as different textures. The masking of the texture 1 as well as texture 2 is applied to create texture dependent divided image. It helped in separating out the background and the outside boundary of image in one texture..

9. Edge Detection: Edges are the limits, line, borders associated with the borders of any object in the image. The edges are taken as the abrupt changes in the shades of the image objects. The sobel filtering are applied first highlight edge borders, and then compute gradient magnitude to provide a picture with one at edges as well as zero for inside areas.

$$f(x) = \frac{I_r - I_l}{2} \left(\text{erf} \left(\frac{x}{\sqrt{2}\sigma} \right) + 1 \right) + I_l$$

I_l and I_r are image intensity at left and right side. σ is the blur scale of image.

10. Edge Erosion: The unwanted images are small lines in the image objects and creates problem in separating of major

segments. Hence, the erosion of such objects is having less than disk size 4 pixel done and performed the modernization. Let E be a Euclidean space or an integer grid, and A a binary image in E . The erosion of the binary image A by the structuring element B is defined by:

$$A \ominus B = \{z \in E | B_z \subseteq A\}$$

where B_z is the translation of B by the vector z , i.e.

$$B_z = \{b + z | b \in B\}, \forall z \in E$$

11. Edge widening: After removal of small and thin boundaries. The leftover boundaries are widened to clearly segment out the important image objects. Implement dilation to objects with a disc size of less than 4 pixels as well as rebuild the image.

12. Thresholding: Thresholding operation is applied to choose segmented bounds of strong luminous intensity on edge elements

13. Watershed Transform: The image consists of different intensities in the gray levels. The higher and lower intensity regions are considered as the maxima and minima location. It is called as peaks and valleys. The watershed transform envisages the distribution of local minima with respect to intensity. It helps in the image objects with separated boundaries.

Let $f \in (D)$ have minima $\{m_k\}_{k \in I}$, for some index set I . The catchment basin (m_i) of a minimum m_i is defined as the set of points $x \in D$ which are topographically closer to m_i than to any other regional minimum :

$$CB(m_i) = \{x \in D | \forall j \neq i: f(m_i) + T_f(x, m_i) < f(m_j) + T_f(x, m_j)\} \quad (i)$$

The watershed of f is the set of points which do not belong to any catchment basin:

$$Watershed = (f) = D \cap \left(\bigcup_{i \in I} CB(m_i) \right)^c$$

Let W be some label W . The watershed transform of f is a mapping $\lambda: D \rightarrow I \cup \{W\}$, such that $\lambda p = i$ if $p \in (m_i)$, and $\lambda p = W$ if $p \in Watershed$.

So, the watershed transform of f assigns labels to the points of D , such that (i) different catchment basins are uniquely labeled, and (ii) a special label W is assigned to all points of the watershed of f .

14. Superimposing of portioned image: All the segmented image objects obtained by entropy features, edge detection and watershed transforms are finally superimposed to account the final segmented copy. Superimpose surface based fragments image over the crisis change connected divided image by alpha mixing.

$$I_{superimposed} = I_1 \cup I_2 \cup I_3$$

15. The segments are termed as the objects representing different kind of features. Extraction of the image features using entropy, wavelet coefficient and textures data for all the images.

16. Development of the training and testing data using the features of segmented image to develop a classification model using the ANFIS algorithm.

17. Application of clustering scheme is applied to generate fuzzy inference system consisting of membership functions and fuzzy rules. The ANN learning scheme is applied to reform the developed fuzzy system using the clustering.

A neuron with label j receiving an input $p_j(t)$ from predecessor neurons consists of the following components: an activation $a_j(t)$, the neuron's state, depending on a discrete time parameter,

an optional threshold θ_j , which stays fixed unless changed by learning,

an activation function that computes the new activation at a given time $t+1$ from $a_j(t)$, θ_j and the net input $p_j(t)$ giving rise to the relation:

$$a_j(t + 1) = f(a_j(t), p_j(t), \theta_j),$$

and an output function computing the output from the activation

$$o_j(t) = f_{out}(a_j(t))$$

The propagation function computes the input to the neuron from the outputs and typically has the form:

$$p_j(t) = \sum_i o_i(t) w_{ij}$$

A bias term can be added, changing the form to the following:

$$p_j(t) = \sum_i o_i(t) w_{ij} + w_{0j}$$

w_{0j} is the bias.

4. Result and Discussion:

We put our word extraction technique from complex damaged photographs to the test using arbitrary images from several online sources. The datasets used to create these photos span a wide range of backgrounds, colours, image quality, font sizes, orientations, and other attributes. The suggested algorithm is put into practise using the MATLAB R2018a platform. On a typical engine running Windows 8.1 Pro 64-bit and an Intel Core i7-4770 CPU clocked at 3.40 GHz, all of the experiments are run. Figures 3 through 26 show some of the complex degraded photos and retrieved texts.



Fig. 2. Original Image

Image with Gaussian noise



Fig. 3. Image with the Gaussian noise

MIDDLEBOROUGH

Fig. 4. Detected Object



Fig 5. Original Image

Image with Gaussian noise

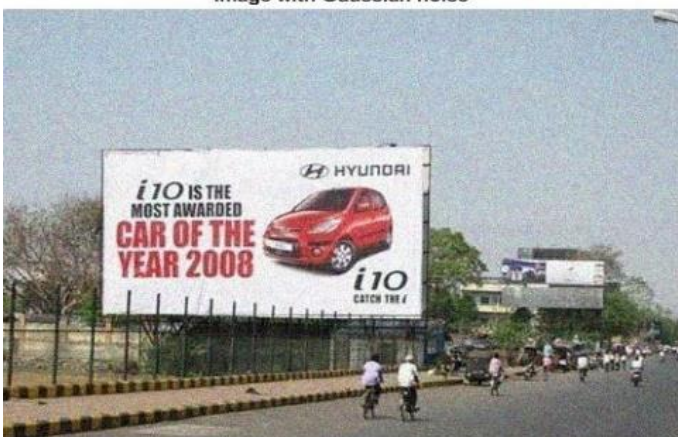


Fig 6. Image with the Gaussian noise

HYUNDAI

i10 IS THE

MOST AWARDED

CAR OF THE

YEAR 2008

i 10

Fig. 7. Detected Object



Fig 8. Original Image

Image with Gaussian noise



Fig 9. Image with the Gaussian noise

INCOME TAX DEPARTMENT
GOVT. OF INDIA
Permanent Account Number Card
ABCDE1234F
APPLICANT NAME
/Fathers Name
APPLICANT'S FAIHER NAME
01/06/1995 Signature

Fig. 10 Detected Object

4. Conclusion:

The performance of text object detection model also dependent on the texture and color intensity in images. The model's scalability to other datasets was previously constrained by the researchers' use of a single hyper-tuned configuration. Our suggested hybrid technique has shown a promising boost in accuracy, especially for text under small object identification. Using the wavelet, it is possible to determine how the watershed and entropy characteristics are configured over the fitted dataset. To increase the model's efficacy, we can additionally employ better morphological methods. The precision rating of the suggested model has enhanced.

References:

- [1] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. I511–I518, Kauai, Hawaii, USA, December 2001.
- [2] R. Lienhart and J. Maydt, "An extended set of Haar-like features for rapid object detection," in Proceedings of the International Conference on Image Processing (ICIP '02), pp. I/900–I/903, Rochester, NY, USA, September 2002.
- [3] P. Viola, J. C. Platt, and C. Zhang, "Multiple Instance boosting for object detection," in Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS '05), vol. 18, pp. 1417–1424, Vancouver, British Columbia, Canada, December 2005.
- [4] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in Proceedings of the 27th IEEE Conference on Computer Vision and Pattern Recognition (CVPR '14), pp. 580–587, Columbus, Ohio, USA, June 2014.
- [5] R. Girshick, "Fast R-CNN," in Proceedings of the 15th IEEE International Conference on Computer Vision (ICCV '15), pp. 1440–1448, Santiago, Chile, December 2015.
- [6] 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*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [8] J. F. Dai, Y. Li, K. M. He et al., "R-FCN: Object Detection via Region-based Fully," in Proceedings of the 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, 2016.
- [9] M. Everingham, L. van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [10] Y. Ren, C. Zhu, and S. Xiao, "Small object detection in optical remote sensing images via modified faster R-CNN," *Applied Sciences*, vol. 8, no. 5, article 813, 2018.
- [11] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 11, pp. 1254–1259, 1998.
- [12] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. S. Torr, and S.-M. Hu, "Global contrast based salient region detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 3, pp. 569–582, 2015.
- [13] C. Chen, M. Y. Liu, O. Tuzel et al., "R-CNN for small object detection," in Asian Conference on Computer Vision, vol. 10115 of Lecture Notes in Computer Science, pp. 214–230, 2016.
- [14] J. S. Lim and W. H. Kim, "Detection of multiple humans using motion information and adaboost algorithm based on Harr-like features," *International Journal of Hybrid Information Technology*, vol. 5, no. 2, pp. 243–248, 2012.
- [15] R. P. Yadav, V. Senthilarasu, K. Kutty, and S. P. Ugale, "Implementation of Robust HOG-SVM based Pedestrian Classification," *International Journal of Computer Applications*, vol. 114, no. 19, pp. 10–16, 2015.
- [16] L. Hou, W. Wan, K.-H. Lee, J.-N. Hwang, G. Okopal, and J. Pitton, "Robust Human Tracking Based on DPM Constrained Multiple-Kernel from a Moving Camera," *Journal of Signal Processing Systems*, vol. 86, no. 1, pp. 27–39, 2017.
- [17] A. Ali and M. A. Bayoumi, "Towards real-time DPM object detector for driver assistance," in Proceedings of the 23rd IEEE International Conference on Image Processing, ICIP 2016, pp. 3842–3846, Phoenix, Ariz, USA, September 2016.
- [18] S. Bell, C. L. Zitnick, K. Bala, and R. Girshick, "Inside-Outside Net: Detecting objects in context with skip pooling and recurrent neural networks," in Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, pp. 2874–2883, Las Vegas, Nev, USA, July 2016.
- [19] T. Kong, A. Yao, Y. Chen, and F. Sun, "HyperNet: towards accurate region proposal generation and joint object detection," in Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 845–853, Las Vegas, Nev, USA, June 2016.
- [20] F. Yang, W. Choi, and Y. Lin, "Exploit all the layers: fast and accurate CNN object detector with scale dependent pooling and cascaded rejection classifiers," in Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, pp. 2129–2137, Las Vegas, Nev, USA, July 2016.