

## Object Detection for Image Captioning

**Kiran Chavan<sup>1\*</sup>, Gayatri Kadam<sup>1</sup>, Ruchika Kankaria<sup>1</sup>, Raksha Kate<sup>1</sup>, Ashvini Ladekar<sup>2</sup>**

<sup>1</sup>B.E. Students, Department of Information Technology, PCCOE, Pune, Maharashtra, India

<sup>2</sup>Professor, Department of Information Technology, PCCOE, Pune, Maharashtra, India

\*Email: [chavankiran38@gmail.com](mailto:chavankiran38@gmail.com)

DOI: <http://doi.org/10.5281/zenodo.2551870>

### Abstract

Generation of description of pictures victimization tongue sentences is gaining a lot of quality of late. It's a difficult task, because it needs not solely understanding a picture, however to translate that visual data into sentence description. So as to caption a picture, we tend to 1st have to be compelled to discover the objects within the image. Object detection has become one amongst the international widespread analysis fields. 1st the paper introduced the distinction between deep learning and machine learning for object detection. Second the techniques for object detection are surveyed. Third it mentioned techniques for object classification.

**Keywords:** Image, Object Detection, Object, Feature extraction, Machine learning

### INTRODUCTION

With the event of mobile internet and also the popularization of varied social media, the number of image knowledge on net has magnified quickly, however groups of people cannot process expeditiously such a big amount of image knowledge. Thus it's expected to hold out these processing mechanically with the help of computer to resolve large-scale visual issues. Object detection technology aims to detect the target objects with the theories and ways of image process and pattern recognition, verify the linguistics classes of those objects, and mark the particular position of the target object within the image. Object detection might be a PC vision strategy for unmistakable items in pictures or recordings. Object detection might be a key yield of profound learning and machine learning calculations. Once humans inspect a photograph or watch a video, we will without delay spot folks, objects, scenes, and visual details. The goal is to show a computer to try and do what comes naturally to humans: to achieve grade of understanding of what a picture contains techniques in machine learning and deep learning have become

widespread approaches to visual perception issues. Each technique learn to spot objects in pictures, however they dissent in their execution.

Using machine learning for visual perception offers the flexibleness to settle on the most effective combination of options and Classifies for learning. It can do correct results with lowest knowledge.

- Benefits of Machine Learning:
- Works higher on small knowledge
- Financially and computationally low-cost
- Easier to interpret

Machine learning techniques are also widespread for visual perception and supply completely different approaches than deep learning:

- HOG feature extraction with associate SVM machine learning model
- Bag of words model with options like SURF and MSER
- The Viola-Jones formula, which might be used to acknowledge a spread of objects, as well as faces and higher bodies.
- Object classification techniques are

Nave Bayes, SVM, k-nearest neighbor and decision tree.

### RELATED WORK

Speeded Up robust features "SURF" formula may be a native feature and descriptor formula that may be utilized in several application like visual perception , SURF use a lot of larger variety of options descriptor from origin image which might scale back contribution of the errors caused by native variation within the average of all feature matching . The strategy for speeded up powerful highlights "SURF" equation might be separated into 3 fundamental advances. begin is "Recognition venture", amid this progression intrigue focuses are chosen at particular areas inside the inception picture, similar to corners, masses and T-intersections and this strategy ought to be heartily. The foremost valuable property of associate interest points it is repeatability. Repeatability specific the dependability of the detector for locating a similar physical interest points below completely different scene conditions. Second step is "Description step", during this step interest points ought to have distinctive identifiers doesn't rely on options scale and rotations that are referred to as descriptor, data of interest points described by descriptor that are vectors that contain information regarding the points itself and also the surroundings. Third step is "Coordinating advance", amid this progression descriptor vectors are analyzed between the thing picture and furthermore the new information or starting point picture, the coordinating score is determined dependent on the space between vectors.

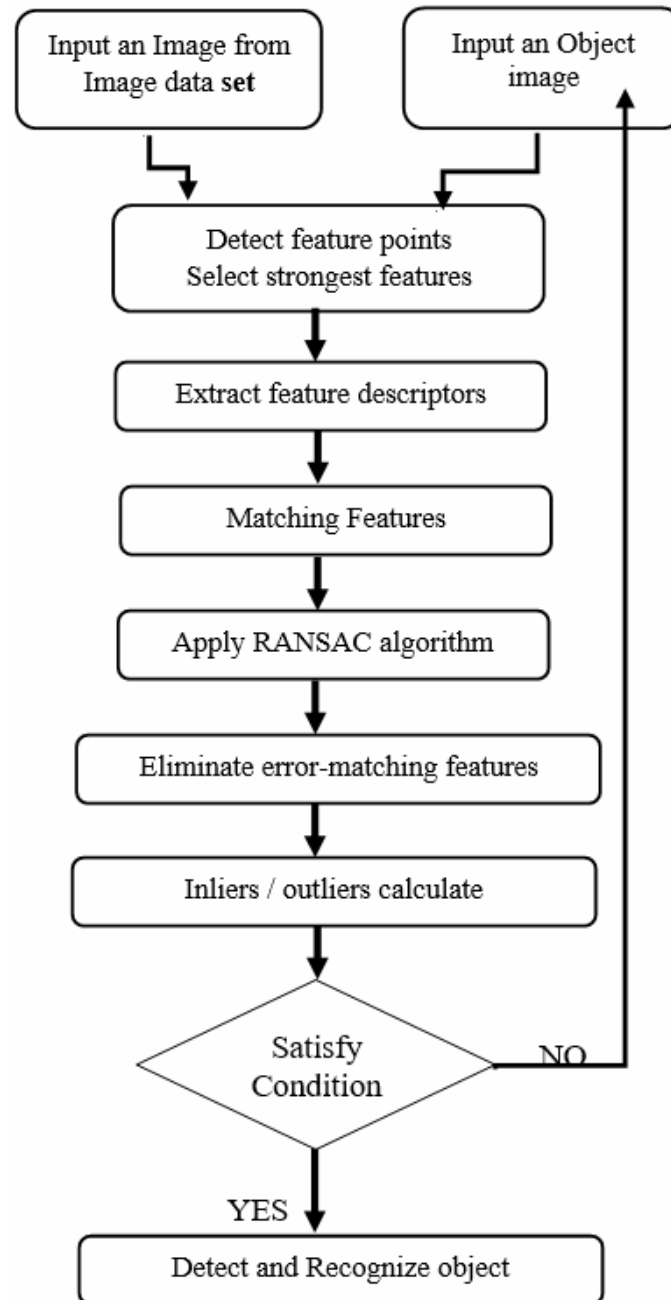
Then if the item is found then provides a message for that and think about the

proportion of matching score and store the end in predefined file, otherwise can provide associate underline message that object wasn't found

### SURF

SURF features descriptor is calculated by the total of Haar wavelet response around interest points. And these may be computed by the conception of integral image. This formula may be utilized in several applications like acknowledge and find of objects, track objects, face recognition, SURF highlights utilized for location the traffic enlist Field Programmable Gate Array (FPGA) that is equipment to strategy video streams progressively.

The interest purpose detection that is described by a vector decision descriptor in SURF formula relies on scale area theory, SURF formula use associate whole number approximations because the determinant of Wellington blob detector which might be computed quick with associate integral image SURF applies completely different sizes from box filters to look and compare interest points, thus box filters has completely different size may be construct the size area and which might be divided into octaves. Scale area illustration is outlined because the convolution of a given image with mathematician kernel. Typically scale area are enforced as a picture pyramid, Scale area may be a continuous operates which might be wont to realize the most values across all attainable scales [9].



*Figure 1:Block Diagram for SURF.*

**HOG**

The histogram of situated inclinations (HOG) is a component descriptor utilized in connection PC vision and picture handling for the article discovery. The procedure includes events of inclination introduction confined parts of an image.

Object detection is required for several embedded vision applications as well as surveillance, advanced driver help systems (ADAS), transportable physics and

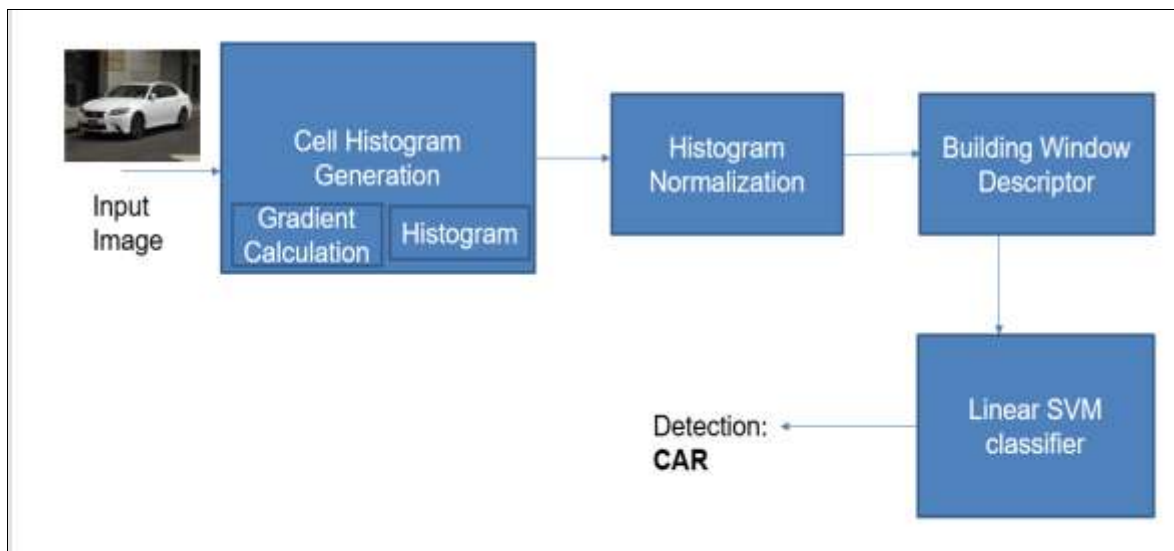
artificial intelligence. For these applications, it's fascinating for object detection to be real-time, strong and energy-efficient. Histogram of oriented Gradients (HOG) may be a wide accepted feature for object detection that provides an affordable trade-off between detection accuracy and complexness compared to various richer options [2].

**HOG feature extraction:** A gradient filter  $[-1 \ 0 \ 1]$  is employed to get a try of

horizontal and vertical gradients for every picture element. The orientation and also the magnitude of the gradient are then calculated from this try, and a histogram of nine bins is generated for the cell. Since the orientation is just wont to select the bar chart bin, the particular angle price of the gradient orientation doesn't have to be compelled to be calculated. Each slope canister might be controlled by correlation vertical and flat inclinations expanded by steady point digressions speaking to receptacles edges. Registering the L2-standard extent of inclinations needs a root activity that is relatively confused for equipment usage. During this work, associate L1-norm is employed for the magnitude to avoid employing a root with no impact to detection accuracy. Every cell needs its neighboring eight cells for the normalization method. Consequently, the ensuing 9-bin cell histogram should be hold on during a column buffer (0.055 Mbit for a single-scale detector), in order that it may be accessed later to reason the normalized bar chart. The standardization is done by separating the 9-canister bar diagram by the square vitality (L2-standard) of everything about four neighboring squares. Dislike the angle extent estimation, utilizing L1-standard for standardization prompts recognition debasement. The root module is

implemented utilizing a non-reestablishing plan and is shared over the four squares. At last, nine continuous mounted reason dividers square measure won't to produce a definitive HOG highlight that might be a 36-measurement vector for each cell [11].

**SVM Classification:** Linear SVM classifiers are typically used for object detection in conjunction with HOG options. During this work, the classifier is trained off-line and also the SVM weights are hold on in associate on-chip SRAM, in order that the detector may be organized for various objects. The bit-width of the SVM loads is diminished to decrease each the memory size and data measure. The 4608 SVM loads are measure to 4-bit marked settled point delineation, with a total memory size of 0.018 Mbit. HOG feature bit-width is chosen to be 9-bit signed fixed-point illustration to take care of the detection accuracy. The HOG feature of every cell is straight away used for classification once it's extracted in order that it's ne'er buffered or recomputed. All computations that grasp the HOG highlight ought to be finished before it's discarded. Thus, the SVM arrangement, that includes a speck item between HOG choices and SVM loads, is done exploitation relate on-the-fly methodology[11].



**Figure 2:**Block Diagram for HOG Object Detection.

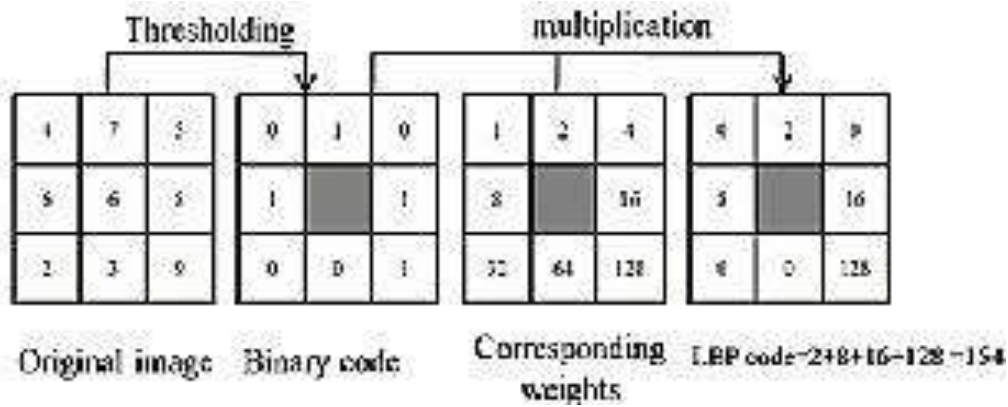
### Local Binary Pattern

The LBP methodology relies on remodeling native binary options of a picture into micro-patterns that may be wont to, as an example moving object detection and face recognition and detection. We tend to show however the LBP feature vectors may be made victimization the quality CNN. Also, we tend to show however straightforward modifications to the quality CNN cell may be won't to create the process of the LBPs simpler. Associate analog readout theme is delineated and also the impact of the analog readout on face recognition accuracy is simulated. The simulations are performed victimization the quality FERET (The facial recognition Technology) information.

In the LBP methodology every picture element is replaced by a binary pattern that comes from the pixel's neighborhood. Every grayscale picture element  $P$  of a picture is employed as a middle of a circle with radius  $r$ . the amount of samples  $M$  determines the number of points that are taken uniformly from the contour of the circle. These points are interpolated from adjacent pixels if required. The sample purposes are compared against the picture element  $P$  one by one with an easy comparison operation which ends up a binary zero if the middle purpose is larger than the present sample point and one otherwise. Once doing this operation as an example dextrorotatory from a definite place to begin the result is a binary pattern with length  $M$ . This operation is illustrated in Figure one. A sensible variety of samples within the read of feature vector length and recognition accuracy has been shown to vary from four to 12. A radius of 1 resembling that every purpose is threshold against its nearest neighbors' was 1st projected for texture recognition, however in [5] the most effective leads to face recognition were earned with a radius of 2. Solely the primary neighborhood is feasible within

the commonest CNN implementations. The effective neighborhood of the LBP will more be magnified by low pass filtering the image in order that rather than taking a sample purpose at a definite picture element its neighborhood is additionally taken into consideration [3]. This is often a linear operation within the CNN.

A histogram is generated from the native binary patterns of a corresponding image. Depending on the application the image may be divided into many spaces in order that every area is represented by a definite histogram. Every distinct pattern is employed as a bin and also the amplitude is that the total of those patterns therein image region. The length of the feature vector of a picture space is thence  $2M$  and if there are  $N$  image areas in a picture the whole concatenated feature vector length is  $N$  times  $2M$ . As an example, with eight samples and twenty image regions the whole feature vector length would be 5120. In a foundation that alittle set of patterns referred to as uniform represent image expeditiously was created. Really the popularity accuracy has shown to be in some cases even higher once victimization this set of patterns than if the patterns with poor discriminative properties were additionally enclosed. The uniform patterns are characterized by the amount of 0-1 and 1-0 transitions in order that second order uniformity is outlined for patterns that have at the most 2 transitions. Second order uniformity seems up to now to be the most effective level of uniformity within the read of recognition potency and have vector length. As an example the patterns 00011000 and 11100111 are second order uniform as a result of the roundness of the LBP operator and on the opposite hand the patterns 10010000 and 00110011 are fourth order uniform. The left and also the right bits square measure thought of as neighbors[8].



**Figure 3: Binary Feature Extraction.**

**Haar Transformation**

Haar remodel is that the simplest associated basic transformation from the area domain to an area frequency domain and may work as an example for orthonormal wavelet transforms steps of Haar wavelet Transform:

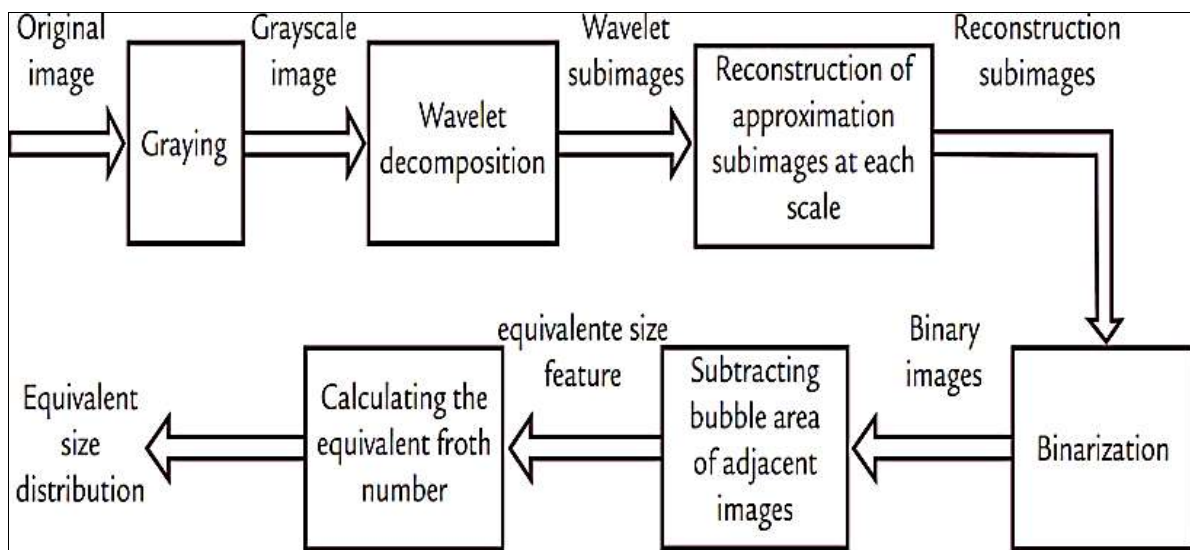
To calculate the Haar remodel of associate array of nsamples:

1. Calculate the average of every pair of samples.(n/2 averages)
2. Calculate the difference between every averageand the samples it had been calculated from. (n/2differences)
3. Now write the 1st half of the array with averages
4. Now write the 2nd half of the array with differences.
5. Repeat the method on the primary

half of the array. Whereas doing this the array size ought to be Power of 2[10].

The main purpose of this methodology is to indicate the impact of the separate Haar wavelet transformation on a picture as Feature Extraction. This is often the combination of a sequence of low-pass and high-pass filters, referred to as a filter bank.

In distinction, wavelets are localized waves. They need their energy focused in time or area and are suited to analysis of transient signals. Whereas Fourier transform and STFT use waves to investigate signals, the wavelet remodel uses wavelets of finite energy[5].



**Figure 4: Haar Wavelet for Object Detection.**

**Table 1: Comparison between Object Detection Techniques.**

| Techniques                      | Advantages   | Disadvantages  |
|---------------------------------|--|--|
| SURF                            | 1. Detect the robust regions properly<br>2. Its accuracy is higher in retrieving the embedded knowledge    | 1. Restricted for native   |
| Haar Wavelet                    | 1. conceptually straightforward.<br>2. It is fast  | 1. no specific use of configuration of visual word positions<br>2. poor at localizing objects among an image |
| Histogram of Oriented Gradients | 1. capable of capturing the pedestrian or object outline/shape better                                      | 1. HOG isn't scale and rotation invariant.   |
| Local Binary Pattern            | 1. its illumination and translation invariant<br>2. efficiently summarizes the native structures of images | 1. They produce rather long histograms, that hamper the recognition speed.                                   |

**SURF:** It detects the robust regions properly and its accuracy is higher in retrieving the embedded knowledge but restricted for native.

**Haar Wavelet:** It is conceptually straightforward and fast but no specific use of configuration of visual word positions and poor at localizing objects among an image

**Histogram of Oriented Gradients:** It is capable of capturing the pedestrian or

object outline/shape better but HOG isn't scale and rotation invariant.

**Local Binary Pattern:** Its illumination and translation invariant and efficiently summarizes the native structures of images but they produce rather long histograms, that hamper the recognition speed.

After studying this we can find out that Histogram of Oriented Gradients is more efficient for object detection

**Table 2: Comparison between Object Classification Techniques**

| Algorithm                        | Features  | Limitations  |
|----------------------------------|---|--|
| K-Nearest neighbor algorithm     | 1. Classes needn't to linearly divisible<br>2. Zero value of the training process | 1. It is sensitive to noisy or irrelevant attributes<br>2. Performance depends on the quantity of dimensions used  |
| Naive bayes Algorithm            | 1. simple to implement<br>2. Great computational potency and classification rate  | 1. The exactness of formula decreases if quantity of information is a smaller amount<br>2. For getting good results it needs a really great deal of data |
| Support Vector Machine Algorithm | 1. High Accuracy<br>2. work well although knowledge isn't linearly separable      | 1. Speed and size demand each within the coaching and testing section<br>2. High complexness and intensive memory demand                                 |

**K-Nearest neighbor algorithm:** It's Classes needn't to linearly divisible and it has zero value of the training process but it is sensitive to noisy or irrelevant attributes.

**Naive bayes Algorithm:** It is simple to implement and it has great computational

potency and classification rate But for getting good results it needs a really great deal of data.

**Support Vector Machine Algorithm:** It has high accuracy and work well although knowledge isn't linearly separable but

speed and size demand each within the coaching and testing section.

After studying this we can find out that Support Vector Machine Algorithm is more efficient for Object Classification

## CONCLUSION

The best approach for visual perception depends on time needed. Sometimes it's very long and it's computationally very costly, and this will be resolved by applying parallelization, tough to configure and not explainable results. The complexness of the hidden layers of deciliter makes it tough to interpret the results or to know the formula mechanism. An object detection model is constructed that consists of accuracy metrics which is able to be more wont to generate captions when detection objects. Using machine learning for visual perception offers the flexibleness to settle on the most effective combination of options and classifiers for learning. It is able to do correct results with lowest knowledge.

## REFERENCES

1. Sagar Nikam. Comparative study of classification techniques in data mining.
2. Master Thesis. *Object Detection and Segmentation Using Contours*. May 2014.
3. Chaudhari Monali, sondur Shanta & Vanjare Gauresh. A review on Face Detection and study of Viola Jones method. July 2015.
4. Wang Zhiqiang & Liu Jun. A review of object detection based on CNN. July 2017.
5. Arora Sangeeta, Yadwinder S. Brar & Kumar Sheo. HAAR Wavelet Transform for Solution of Image Retrieval. 2014.
6. Aghasi S. Poghosyan. Image Caption Generation and Object Detection via a Single Model. 2017.
7. Singh Sisodiya, Deep learning

- approaches for image caption generation Ayushman. 2016.
8. O.Lahdenoja, Laiho M. & Paasio A. Local binary pattern feature vector extraction with CNN. *IEEE*. 2005.
9. Herbert Bay, Andreas Ess, Tinne Tuytelaars, & Luc Van Gool. Speeded-Up Robust Features (SURF). *Speeded-Up Robust Features (SURF)*.
10. Shantikumar Singh Y., Pushpa Devi B. & Manglem Singh Kh. Image Compression with Haar Wavelet Transform. *International Journal of Computer Applications*. 2015.
11. Amr Suleiman & Vivienne Sze. Energy-Efficient HOG-based Object Detection at 1080HD 60 fps with Multi-Scale Support. *IEEE International Workshop on Signal Processing Systems (SiPS)*. 2014.
12. Papageorgiou C.P. A general framework for object detection. *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*.
13. Felzenszwalb Pedro F., Ross B. Girshick, David McAllester & Deva Ramanan. Object Detection with Discriminatively Trained Part-Based Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. September 2010; 32(9).
14. Gabriella Csurka, Christopher R. Dance, Lixin Fan, Jutta Willamowski & Cédric Bray. Visual Categorization with Bags of Keypoints. *Xerox Research Centre Europe 6, chemin de Maupertuis 38240 Meylan, France*.

**Cite this article as:** Kiran Chavan, Gayatri Kadam, Ruchika Kankaria, Raksha Kate, & Ashvini Ladekar. (2019). Object Detection for Image Captioning. *Journal of Image Processing and Artificial Intelligence*, 5(1), 17–24. <http://doi.org/10.5281/zenodo.2551870>