

HoughNet: Integrating near and long-range evidence for visual detection

Nermin Samet[†], Samet Hicsonmez, and Emre Akbas

Abstract—This paper presents HoughNet, a one-stage, anchor-free, voting-based, bottom-up object detection method. Inspired by the Generalized Hough Transform, HoughNet determines the presence of an object at a certain location by the sum of the votes cast on that location. Votes are collected from both near and long-distance locations based on a log-polar vote field. Thanks to this voting mechanism, HoughNet is able to integrate both near and long-range, class-conditional evidence for visual recognition, thereby generalizing and enhancing current object detection methodology, which typically relies on only local evidence. On the COCO dataset, HoughNet’s best model achieves 46.4 AP (and 65.1 AP₅₀), performing on par with the state-of-the-art in bottom-up object detection and outperforming most major one-stage and two-stage methods. We further validate the effectiveness of our proposal in other visual detection tasks, namely, video object detection, instance segmentation, 3D object detection and keypoint detection for human pose estimation, and an additional “labels to photo” image generation task, where the integration of our voting module consistently improves performance in all cases. Code is available at <https://github.com/nerminsamet/houghnet>.

Index Terms—Object detection, voting, bottom-up recognition, Hough Transform, video object detection, instance segmentation, 3D object detection, human pose estimation, image-to-image translation, label-to-image translation.



1 INTRODUCTION

Deep learning has brought on remarkable improvements in object detection. Performance on widely used benchmark datasets, as measured by mean average-precision (mAP), has at least doubled (from 0.33 mAP [1] [2] to 0.80 mAP on PASCAL VOC [3]; and from 0.2 mAP [4] to around 0.5 mAP on COCO [5]) in comparison to the pre-deep-learning, shallow methods. Current state-of-the-art, deep learning based object detectors [5], [6], [7], [8] predominantly follow a top-down approach where objects are detected holistically via rectangular region classification. This was not the case with the pre-deep-learning methods. The bottom-up approach was a major research focus as exemplified by the prominent voting-based (the Implicit Shape Model [9]) and part-based (the Deformable Parts Model [10]) methods. However, today, among deep learning based object detectors, the bottom-up approach has not been sufficiently explored with a few exceptions (e.g. CornerNet [11], ExtremeNet [12]).

In this paper, we propose HoughNet, a one-stage, anchor-free, voting-based, bottom-up object detection method. HoughNet is based on the idea of voting, inspired by the Generalized Hough Transform [13], [14]. In its most generic form, the goal of GHT is to detect a whole shape based on its parts. Each part produces a hypothesis, i.e. casts its vote, regarding the location of the whole shape. Then, the location with the most votes is selected as the result. Similarly, in HoughNet, the presence of an object belonging to a certain class at a particular location is determined by the sum of the class-conditional votes cast on that location (Fig. 1). HoughNet processes the input image using a con-

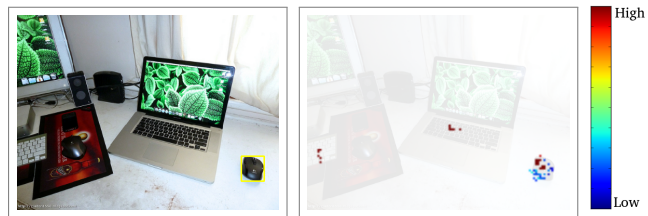


Fig. 1: (Left) A sample “mouse” detection, shown with yellow bounding box, by HoughNet. (Right) The locations that vote for this detection. Colors indicate vote strength. In addition to the local votes originating from the mouse itself, there are strong votes from nearby “keyboard” objects, which shows that HoughNet is able to utilize both short and long-range evidence for detection. More examples can be seen in Fig. 7.

volitional neural network to produce an intermediate score map per class. Scores in these maps indicate the presence of visual structures that would support the detection of an object instance. These structures could be object parts, partial objects or patterns belonging to the same or other classes. We name these score maps as “visual evidence” maps. Each spatial location in a visual evidence map votes for target areas that are likely to contain objects. Target areas are determined by placing a log-polar grid, which we call the “vote field,” centered at the voter location. The purpose of using a log-polar vote field is to reduce the spatial precision of the vote as the distance between voter location and target area increases. This is inspired by foveated vision systems found in nature, where the spatial resolution rapidly decreases from the fovea towards the periphery [15]. Once all visual evidence is processed through voting, the accumulated votes are recorded in object presence maps,

N. Samet and E. Akbas are at the Dept. of Computer Engineering, Middle East Technical University (METU), Ankara, Turkey. S. Hicsonmez is at the Dept. of Computer Engineering, Hacettepe University, Ankara, Turkey. E-mail: nermin@ceng.metu.edu.tr, samethicsonmez@hacettepe.edu.tr, emre@ceng.metu.edu.tr

[†] Corresponding author.

where the peaks indicate the presence of object instances.

Current state-of-the-art object detectors rely on local or short-range visual evidence to decide whether there is an object at that location (as in top-down methods) or an important keypoint such as a corner (as in bottom-up methods). On the other hand, HoughNet is able to integrate both short and long-range visual evidence at the same time through voting. An example is illustrated in Fig. 1, where the detected mouse gets strong votes from two keyboards, one of which is literally at the other side of the image. In another example (Fig. 7, row 2, col 1), a ball on the right-edge of the image is voting for the baseball bat on the left-edge. On the COCO dataset, HoughNet achieves comparable results with the state-of-the-art bottom-up detector CenterNet [16], while being the fastest among bottom-up detectors. It outperforms prominent one-stage (RetinaNet [5]) and two-stage detectors (Faster R-CNN [6], Mask R-CNN [17]). To further show the effectiveness of our approach, we used the voting module of HoughNet in five other tasks, namely, video object detection, instance segmentation, 3D object detection, human pose estimation and “labels to photo” image generation, and showed that the performances are consistently improved in all cases.

To deal with the large number of training experiments in model development and ablation studies, we curated “COCO minitrain”, a mini training set for COCO. We validated COCO minitrain in two ways by (i) showing that the COCO val2017 performance of a model trained on COCO minitrain is strongly positively correlated with the performance of the same model trained on COCO train2017, and (ii) showing that COCO minitrain set preserves the object instance statistics.

This paper extends our previous work [18] in several ways. (i) The original HoughNet [18] applies voting only in the spatial domain (for object detection in still images). We extended this idea to the temporal domain by developing a new method, which takes the difference of features from two frames, and applies spatial and temporal voting using our “temporal voting module” to detect objects. We showed the effectiveness of this method in video object detection (Sections 3.4 and 4.3.1). (ii) We showed that our Hough voting module is effective not only in object detection but also in other visual detection tasks as well: instance segmentation (Section 4.3.2), human pose estimation, i.e. keypoint detection (Section 4.3.4) and 3D object detection (Section 4.3.3). (iii) We added an analysis section (Section 4.2.4) where we examined the relations between vote-giver and vote-getter classes, and evaluated HoughNet’s performance using localisation, recall, precision (LRP) metrics [19], [20] to identify and compare the source of errors. (iv) We developed a “scalable” variant of HoughNet where the number of voting operations does not depend on the number of object classes (Section 4.2.5) and showed that it dramatically improves inference speed with a slight drop in accuracy. (v) And finally, we improved the source code of HoughNet¹ by increasing its modularity and training speed.

Our main contribution in this work is HoughNet, a voting-based detection method that is able to integrate near and long-range evidence for visual detection tasks.

2 RELATED WORK

Methods using log-polar fields/representations. Many biological systems have foveated vision where the spatial resolution decreases from the fovea (point of fixation) towards the periphery. Inspired by this phenomenon, computer vision researchers have used log-polar fields for many different purposes including shape description [21], feature extraction [22] and foveated sampling/imaging [23].

Non-deep, voting-based object detection methods. In the pre-deep learning era, generalized Hough Transform (GHT) based voting methods have been used for object detection. The most influential work was the Implicit Shape Model (ISM) [9]. In ISM, Leibe *et al.* [9] applied GHT for object detection/recognition and segmentation. During the training of the ISM, first, interest points are extracted and then a visual codebook (i.e. dictionary) is created using an unsupervised clustering algorithm applied on the patches extracted around interest points. Next, the algorithm matches the patches around each interest point to the visual word with the smallest distance. In the last step, the positions of the patches relative to the center of the object are associated with the corresponding visual words and stored in a table. During inference, patches extracted around interest points are matched to closest visual words. Each matched visual word casts votes for the object center. In the last stage, the location that has the most votes is identified, and object detection is performed using the patches that vote for this location. Later, ISM was further extended with discriminative frameworks [24], [25], [26], [27], [28]. Okada [24] ensembled randomized trees using image patches as voting elements. Similarly, Gall and Lempitsky [25] proposed to learn a mapping between image patches and votes using random forests. In order to fix the accumulation of inconsistent votes of ISM, Razavi *et al.* [26] augmented the Hough space with latent variables to enforce consistency between votes. In Max-margin Hough Transform [27], Maji and Malik showed the importance of learning visual words in a discriminative max-margin framework. Barinova *et al.* [28] detected multiple objects using energy optimization instead of non-maxima suppression peak selection of ISM.

HoughNet is similar to ISM and its variants described above only at the idea level as all are voting based methods. There are two major differences: (i) HoughNet uses deep neural networks for part/feature (i.e. visual evidence) estimation, whereas ISM uses hand-crafted features; (ii) ISM uses a discrete set of visual words (obtained by unsupervised clustering) and each word’s vote is exactly known (stored in a table) after training. In HoughNet, however, there is not a discrete set of words and vote is carried through a log-polar vote field which takes into account the location precision as a function of target area.

Bottom-up object detection methods. Apart from the classical one-stage [5], [7], [29], [30], [31] vs. two-stage [6], [17] categorization of object detectors, we can also categorize the current approaches into two: top-down and bottom-up. In the top-down approach [5], [6], [7], [29], a near-exhaustive list of object hypotheses in the form of rectangular boxes are generated and objects are predicted in a holistic manner based on these boxes. Designing the hypotheses space (e.g.

1. <https://github.com/nerminsamet/houghnet>

parameters of anchor boxes) is a problem by itself [32]. Typically, a single template is responsible for the detection of the whole object. In this sense, recent anchor-free methods [33], [34] are also top-down. On the other hand, in the bottom-up approach, objects *emerge* from the detection of parts or sub-object structures. For example, in CornerNet [11], top-left and bottom-right corners of objects are detected first, and then, they are paired to form whole objects. Following CornerNet, ExtremeNet [12] groups extreme points (e.g. left-most, etc.) and center points to form objects. Together with corner pairs of CornerNet [11], CenterNet [16] adds center point to model each object as a triplet. HoughNet follows the bottom-up approach based on a voting strategy: object presence score is voted (aggregated) from a wide area covering short and long-range evidence.

Deep, voting-based object detection methods. Qi et al. [35] apply Hough voting for 3D object detection in point clouds. Sheshkus et al. [36] utilize Hough transform for vanishing points detection in the documents. For automatic pedestrian and car detection, Gabriel et al. [37] proposed using discriminative generalized Hough transform for proposal generation in edge images, later to further refine the boxes, they fed these proposals to deep networks. In the deep learning era, we are not the first to use a log-polar vote field in a voting-based model. Lifshitz et al. [38] used a log-polar map to estimate keypoints for single person human pose estimation. Apart from the fact that they are tackling a different problem (human pose estimation), there are several subtle differences. First, they prepare ground truth voting maps for each keypoint such that keypoints vote for every other one depending on its relative position in the log polar map. This requires manually creating static voting maps. Specifically, their model learns $H \times W \times R \times C$ voting map, where R is the number of bins and C is the augmented keypoints. In order to produce keypoint heatmaps, they perform vote aggregation at test phase. Second, this design restricts the model to learn only the keypoint locations as voters. When we consider the object detection task and its complexity, it is not trivial to decide on vote-giving locations for objects or prepare ground-truth voting maps as in human pose estimation. Moreover, this design limits the voters to reside only inside of the object (e.g. person), whereas in our approach an object could get votes from far away regions. To overcome these issues, unlike their model we apply vote aggregation during training (they perform vote aggregation only at test phase). This allows us to expose the latent patterns between objects and voters for each class. In this way, our voting module is able to get votes from non-labeled objects (e.g. “candle”, which is not a COCO class, voting for dining table; see the last row of Fig. 7). To the best of our knowledge, we are the first to use a log-polar vote field in a voting-based deep learning model to integrate the **long range interactions** for object detection.

Similar to HoughNet, Non-local neural networks (NLNN) [39] and Relation networks (RN) [40] integrate long-range features. As a fundamental difference, in NLNN, the relative displacement between interacting features is not taken into account. However, HoughNet uses this information encoded through the regions of the log-polar vote field. RN models object-object relations explicitly for proposal-based two-stage detectors.

Concurrently with HoughNet [18], several object detectors have been introduced [41], [42], [43]. Among them, the transformer-based DETR [41], in particular, shares with HoughNet the idea of using both short and long range interactions. In DETR, features at a specific location are updated through interactions with features at other locations using encoder-decoder based transformers [44]. Although locations of features are taken into account using positional encoding, DETR models all possible pairs of features. Unlike DETR, HoughNet explicitly takes into account the spatial precision of the vote depending on the relative displacement between voter location and target (voted) area. Although attention-based neural networks and HoughNet share the idea of using long-range interactions, an in-depth comparison is beyond the scope of this work. Briefly; there are a variety of ways to encode location in attention-based neural networks: (i) simply no encoding – location, whether absolute or relative, is ignored [39], (ii) using sine and cosine functions with varying frequencies [41], [44], (iii) using a general function [45]. It is not trivial to compare these encodings to the explicit encoding of relative location using a log-polar field in HoughNet. Another fundamental difference is that transformers operate at the feature level by modifying features through interactions with other features, whereas HoughNet operates at the detection score level.

3 HOUGHNET: MODELS AND METHOD

Brief overview. In HoughNet, the input image first passes through a backbone convolutional neural network (CNN), the output of which is connected to three different branches carrying out the predictions of (i) visual evidence scores, (ii) objects’ bounding box dimensions (width and height), and (iii) objects’ center location offsets. The first branch is where the voting occurs. Before we describe our voting mechanism in detail, we first introduce the log-polar vote field below. The overall processing pipeline of HoughNet is illustrated in Fig. 2.

3.1 The Log-polar “Vote Field”

We use the set of regions in a standard log-polar coordinate system to define the regions through which votes are collected. A log-polar coordinate system is defined by the number and radii of eccentricity bins (or rings) and the number of angle bins. We call the set of cells or regions formed in such a coordinate system as the “vote field” (Fig. 3). In our experiments, we used different vote fields with different parameters (number of angle bins, etc.) as explained in the Experiments section. In the following, R denotes the number of regions in the vote field and K_r is the number of pixels in a particular region r . $\Delta_r(i)$ denotes the relative spatial coordinates of the i^{th} pixel in the r^{th} region, with respect to the center of the field. We implement the vote field as a fixed-weight (non-learnable) transposed-convolution filter as further explained below.

3.2 Voting Module

After the input image is passed through the backbone network and the “visual evidence” branch, the voting module of HoughNet receives C tensors $\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_C$, each of size

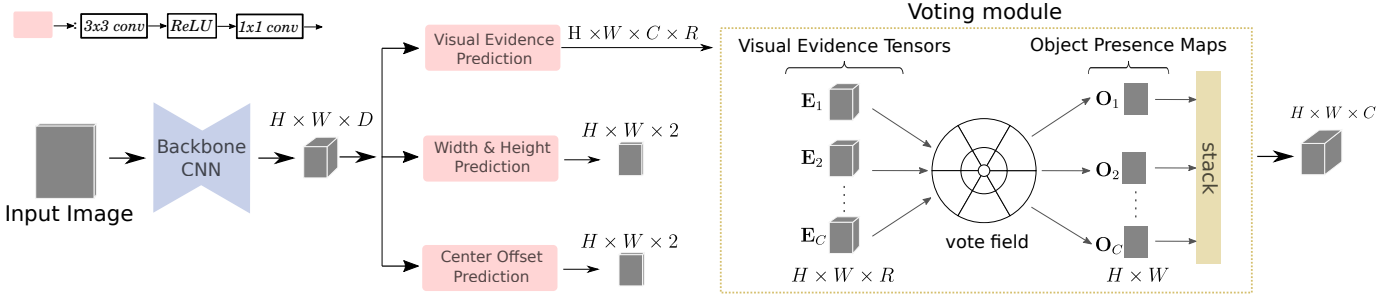


Fig. 2: Overview of the processing pipeline of HoughNet.

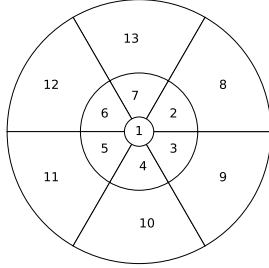


Fig. 3: A log-polar “vote field” used in the voting module of HoughNet. Numbers indicate region ids. A vote field is parametrized by the number of angle bins, and the number and radii of eccentricity bins, or rings. In this particular vote field, there are a total of 13 regions, 6 angle bins and 3 rings. The radii of the rings are 2, 8 and 16, respectively.

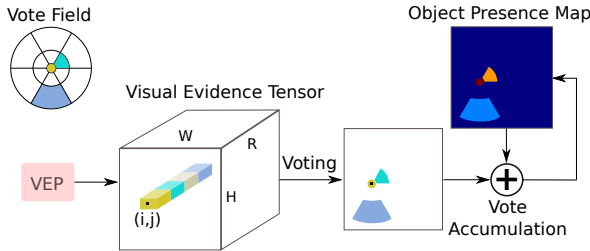


Fig. 4: Voting process for a single class. Visual evidence prediction (VEP) branch outputs $H \times W \times R$ dimensional visual evidence tensor for the class. Here the voting process is illustrated for just a single location (i, j) and for 3 regions of the vote field shown with yellow, green and blue colors. The values at (i, j) corresponding to these 3 regions are added as votes to the appropriate locations in the Object Presence Map tensor. These locations are determined by the vote field shown at top left. Votes from all locations are similarly accumulated in the Object Presence Map.

$H \times W \times R$, where C is the number of classes, H and W are spatial dimensions and R is the number of regions in the vote field. Each of these tensors contains class-conditional (i.e. for a specific class) “visual evidence” scores. The job of the voting module is to produce C “object presence” maps O_1, O_2, \dots, O_C , each of size $H \times W$. Then, peaks in these maps will indicate the presence of object instances. The voting process, which converts the visual evidence tensors (e.g. E_c) to object presence maps (e.g. O_c), works as described below.

The voting process is best explained using an example. Suppose we are to process the visual evidence at the i^{th} row, j^{th} column and the r^{th} channel of a visual evidence tensor E . We first place our vote field (Fig. 3) centered at (i, j) on the r^{th} channel, which is a 2D map. The region r of the vote field marks the target area to be voted on, whose coordinates can be calculated by adding the coordinate offsets $\Delta_r(\cdot)$ to (i, j) . Then, we add the visual evidence score $E(i, j, r)$ to the target area of the object presence map. Note that this operation can be efficiently implemented using the “transposed convolution” (or “deconvolution”) operation. Visual evidence scores from locations other than (i, j) are processed in the same way and the scores are accumulated in the object presence map. We formally define this procedure in Algorithm 1, which takes in a visual evidence tensor as input and produces an object presence map². Fig. 4 also illustrates the defined voting process on a single class.

Algorithm 1 Voting process.

Input: Visual evidence tensor E_c , Vote field relative coordinates Δ
Output: Object presence map O_c
 Initialize O_c with all zeros
for each pixel (i, j, r) in E_c **do**
 /* K_r : number of pixels in the vote field region r */
 for $k = 1$ to K_r **do**
 $(y, x) \leftarrow (i, j) + \Delta_r(k)$
 $O_c(y, x) \leftarrow O_c(y, x) + \frac{1}{K_r} E_c(i, j, r)$
 end for
end for

3.3 Network Architecture

Our network architecture design follows that of “Objects as Points” (OAP) [34]. HoughNet consists of a backbone and three subsequent branches which predict (i) visual evidence scores, (ii) bounding box widths and heights, and (iii) center offsets. Our voting module is attached to the visual evidence branch (Fig. 2).

The output of the backbone network is a feature map of size $H \times W \times D$, which is the result of an input image of size $4H \times 4W \times 3$. The backbone’s output is fed to all three branches. Each branch has one convolutional layer with 3×3 filters followed by a ReLU layer and another

2. We provide a step-by-step animation of the voting process in the supplementary material and at <https://shorturl.at/iIOP2>.

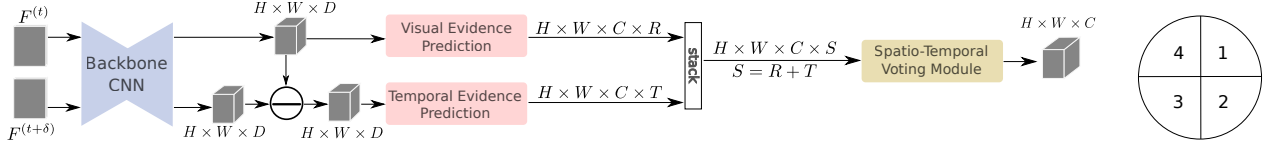


Fig. 5: (Left) Overall processing pipeline of HoughNet for video object detection. (Right) Temporal vote field.

convolutional layer with 1×1 filters. The visual evidence branch outputs $H \times W \times C \times R$ sized output where C and R correspond to the number of classes and vote field regions, respectively. The width & height prediction branch outputs $H \times W \times 2$ sized output which predicts heights and widths for each object center. Finally, the center offset branch predicts relative displacement of center locations across the spatial axes. These offsets help recover the lost precision of the center points due to down-sampling operations through the network. Both the width & height branch and the center offset branch are class-agnostic.

Objective functions. For the optimization of the visual evidence branch, we use the modified focal loss [5] introduced in CornerNet [11] (also used in [12], [34]). We optimize the center offset branch using the L_1 loss as the other bottom-up detectors [11], [12], [34] do. Finally, for the width & height prediction branch, we use L_1 loss by scaling the loss by 0.1 as proposed in OAP [34]. The overall loss is the sum of the losses from all branches.

3.4 Spatio-temporal Voting

Unlike static images, videos have rich temporal information. In order to benefit from the temporal clues in videos, researchers developed several methods to aggregate information locally and globally using two or more frames [46], [47], [48], [49]. Similarly, we extend HoughNet with a new temporal voting module to incorporate temporal information using an additional (auxiliary) frame.

Before describing the temporal voting process, we first explain the temporal log-polar vote field. The temporal vote field has 4 regions, 90° angle bin and a single ring with radii 8. Each region of the temporal vote field stands for a motion direction. For example, the temporal vote-field region with id 1 in Fig. 5, corresponds to relative motion in $+x$ and $+y$ direction. When the temporal voting field is centered at a voter location, it votes for the locations in the target area depending on relative motion direction between frames.

We show the general processing pipeline for video object detection in Fig. 5. In addition to the reference frame at time t , we choose a random auxiliary frame within δ time interval around the reference frame. First, both frames pass through the backbone and we obtain feature maps of size $H \times W \times D$ for each. Then, we calculate motion features between these frames by subtracting the feature map of the auxiliary frame from the feature map of the reference frame. Bertasius et al. showed that using motion features is effective for pose detection in videos [50]. They interpret these features as discriminative motion features, where the network learns to ignore uninformative motion regions while focusing on discriminative motion cues.

Later, relative motion features are fed to the temporal evidence branch. Finally, visual and temporal evidence tensors

are stacked and forwarded to the spatio-temporal voting module to output $H \times W \times C$ dimensional object presence maps. Spatio-temporal voting aggregates votes in the same way as described in Section 3.2. Temporal evidence branch has the same architecture as the visual evidence branch.

4 EXPERIMENTS

In this section, we show the effectiveness of our proposed method on object detection and other visual detection tasks, namely, video object detection, instance segmentation, 3D object detection and keypoint detection (2D human pose estimation). Since we mostly focus on object detection, the object detection section is more detailed compared to others. It includes (i) ablation experiments through which we studied how different parameters of the vote field affect the performance, (ii) comparison with baseline, (iii) comparison with state-of-the-art, (iii) an analysis section and (iv) a scalable voting approach where the number of “visual evidence tensors” does not depend on the number of object classes. The sections for other tasks mostly contain comparisons with baseline. We also include an image generation task in the context of “label-to-photo” translation at the end. In order to handle a large volume of ablation experiments, we created a small training set for COCO, called “COCO minitrain”, which is described first in the following.

We ran our experiments on 4 V100 GPUs. For training, we used 512×512 images unless stated otherwise. The training setup is not uniform across different experiments, mainly due to different backbones. However, the inference pipeline is common for all HoughNet models. We extract center locations by applying a 3×3 max pooling operation on object presence heatmaps and pick the highest scoring 100 points as detections. Then, we adjust these points using the predicted center offset values. Final bounding boxes are generated using the predicted width & height values on these detections.

4.1 COCO minitrain

To facilitate model development and faster analysis in ablation experiments, we carefully curated a mini training set for COCO, dubbed “COCO minitrain”. It is a subset of the COCO train2017 dataset, containing 25K images (about 20% of train2017) and around 184K objects across 80 object categories. We randomly sampled these images from the full set while preserving the following three quantities as much as possible: (i) proportion of object instances from each class, (ii) overall ratios of small, medium and large objects, (iii) per class ratios of small, medium and large objects.

To validate COCO minitrain, we computed the correlation between the val2017 performance of a model

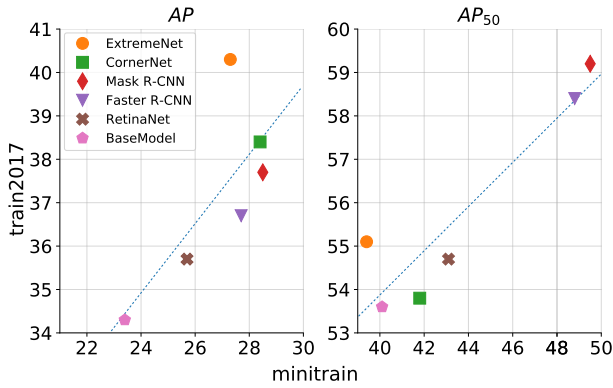


Fig. 6: Performance validation for COCO *minitrain*. The x-axis is the *val2017* performance (AP on the left, AP50 on the right) of a model when it is trained on *minitrain*. The y-axis is the performance of the model when it is trained on the full COCO training set, *train2017*. Fitted lines with Pearson correlation coefficients 0.74 and 0.92, respectively for AP and AP₅₀, show strong positive correlation.

when it is trained on *minitrain* with the same as when it is trained on *train2017*. Over six different object detectors (Faster R-CNN, Mask R-CNN, RetinaNet, CornerNet, ExtremeNet and HoughNet), the Pearson correlation coefficients turned out to be 0.74 and 0.92 for AP and AP₅₀, respectively (Fig. 6), which indicate strong positive correlation. Further details on *minitrain* and the dataset itself can be found at <https://github.com/giddyup/coco-minitrain>.

4.2 Hough Voting for Object Detection

In the evaluation of object detection models, we follow the other bottom-up methods [11], [12], [34] and use two modes: (i) single-scale, horizontal-flip testing (SS testing mode), and (ii) multi-scale, horizontal-flip testing (MS testing mode). In MS, we use the following scale values, 0.6, 1.0, 1.2, 1.5, 1.8. To merge augmented test results, we use Soft-NMS [51], and keep the top 100 detections. All tests are performed on a single V100 GPU.

4.2.1 Ablation Experiments

Here we analyze the effects of the number of angle and ring bins of the vote field on performance. Models are trained on COCO *minitrain* and evaluated on *val2017* set with SS testing mode. The backbone is Resnet-101 [3]. In order to get higher resolution feature maps, we add three deconvolution layers on top of the default Resnet-101 network, similar to [52]. We add 3 × 3 convolution filters before each 4 × 4 deconvolution layer, and put batchnorm and ReLU layers after convolution and deconvolution filters. We trained the network with a batch size of 44 for 140 epochs with Adam optimizer [53]. Initial learning rate 1.75 × 10⁻⁴ was divided by 10 at epochs 90 and 120.

Angle bins. We started with a large, 65 by 65, vote field with 5 rings. We set the radius of these rings from the most inner one to the most outer one as 2, 8, 16, 32 and 64 pixels, respectively. We experimented with 60°, 90°, 180° and 360° bins. We do not split the center ring (i.e. region with id

TABLE 1: Ablation experiments for the vote field. (a) Effect of angle bins on performance. Vote field with 90° has the best performance (considering AP and AP₅₀). (b) Effect of central and peripheral regions. Here, the angle bin is 90° and the ring count is four. Disabling any of center or periphery hurts performance, cf. (a). (c) Effect of number of rings. Angle is 90° and vote field size is updated according to the radius of the last ring. Using 3 rings yields the best result. It is also the fastest model.

Model	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	FPS
60°	24.6	41.3	25.0	8.2	27.7	36.2	3.4
90°	24.6	41.5	25.0	8.2	27.7	36.2	3.5
180°	24.5	41.1	24.8	8.1	27.7	36.3	3.5
360°	24.6	41.1	25.1	8.0	27.8	36.3	3.5
(a) Varying the Number of Angle Bins							
Only Center	23.8	39.5	24.5	7.9	26.8	34.7	3.5
No Center	24.4	40.9	24.9	7.4	27.6	37.1	3.3
Only Context	23.6	39.7	24.2	7.4	26.4	35.9	3.4
(b) Effectiveness of Votes from Center or Periphery							
5 Rings	24.6	41.5	25.0	8.2	27.7	36.2	3.5
4 Rings	24.5	41.1	25.3	8.2	27.8	36.1	7.8
3 Rings	24.8	41.3	25.6	8.4	27.6	37.5	15.6
(c) Varying Ring Counts							

TABLE 2: Comparing our voting module to an equivalent dilated convolution filter, in terms of number of parameters and the spatial filter size, on COCO *val2017* set. Models are trained on COCO *train2017* and results are presented on SS testing mode.

Method	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Baseline	36.2	54.8	38.7	16.3	41.6	52.3
+ Dilated Conv.	36.6	56.1	39.2	16.7	42.0	53.6
+ Voting Module	37.3	56.6	39.9	16.8	42.6	55.2

1 in Fig. 3) into further regions. Results are presented in Table 1a. For the 180° experiment, we divide the vote field horizontally. 90° yields the best performance considering both AP and AP₅₀. We used this setting in the rest of the experiments.

Effects of center and periphery. We conducted experiments to analyze the importance of votes coming from different rings of the vote field. Results are presented in Table 1b. In the *Only Center* case, we only keep the center ring and disable the rest. In this way, we only aggregate votes from features of the object center directly, which corresponds to a traditional object detector where only local (short-range) evidence is used. This experiment shows that votes from outer rings help improve performance. For the *No Center* case, we only disable the center ring. We observe that there is only 0.2 decrease in AP. This suggests that the evidence for successful detection is embedded mostly around the object center not directly inside the object center. In order to observe the power of long-range votes, we conducted another experiment called “Only Context,” where we disabled the two most inner rings and used only the three outer rings for vote aggregation. This model reduced AP by 1.0 point compared to the full model.

Ring count. To find out how far an object should get votes from, we discard outer ring layers one by one as

TABLE 3: HoughNet results on COCO `val2017` set for different training setups. † indicates initialization with CornerNet weights, * indicates initialization with ExtremeNet weights. Results are given for SS and MS testing modes, respectively.

Models	Backbone	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L	FPS
Base	R-101	36.0 / 40.7	55.2 / 60.6	38.4 / 43.9	16.2 / 22.5	41.7 / 44.2	52.0 / 55.7	3.5 / 0.5
Base	R-101-DCN	37.3 / 41.6	56.6 / 61.2	39.9 / 44.9	16.8 / 22.6	42.6 / 44.8	55.2 / 58.8	3.3 / 0.4
Light	R-101-DCN	37.2 / 41.5	56.5 / 61.5	39.6 / 44.5	16.8 / 22.5	42.5 / 44.8	54.9 / 58.4	14.3 / 2.1
Light	HG-104	40.9 / 43.7	59.2 / 61.9	44.1 / 47.3	23.8 / 27.5	45.3 / 45.9	52.6 / 56.2	6.1 / 0.8
Light	HG-104†	41.7 / 44.7	60.5 / 63.2	45.6 / 48.9	23.9 / 28.0	45.7 / 47.0	54.6 / 58.1	5.9 / 0.8
Light	HG-104*	43.0 / 46.1	62.2 / 64.6	46.9 / 50.3	25.5 / 30.0	47.6 / 48.8	55.8 / 59.7	5.7 / 0.8

TABLE 4: Comparison of object detection with baseline (OAP) on `val2017`. Results are given for SS and MS test modes, respectively.

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L	$moLRP \downarrow$
Baseline <i>w</i> R-101-DCN	36.2 / 39.2	54.8 / 58.6	38.7 / 41.9	16.3 / 20.5	41.6 / 42.6	52.3 / 56.2	71.1 / 68.3
+ Voting Module	37.2 / 41.5	56.5 / 61.5	39.6 / 44.5	16.8 / 22.5	42.5 / 44.8	54.9 / 58.4	69.9 / 66.6
Baseline <i>w</i> HG-104	42.2 / 45.1	61.1 / 63.5	46.0 / 49.3	25.2 / 27.8	46.4 / 47.7	55.2 / 60.3	66.1 / 63.9
+ Voting Module	43.0 / 46.1	62.2 / 64.6	46.9 / 50.3	25.5 / 30.0	47.6 / 48.8	55.8 / 59.7	65.6 / 63.1

presented in Table 1c. The models with 5 rings, 4 rings and 3 rings have 17, 13 and 9 voting regions and 65, 33 and 17 vote field sizes, respectively. The model with 3 rings yields the best performance on AP metric and is the fastest one at the same time. On the other hand, the model with 5 rings yields 0.2 AP_{50} improvement over the model with 3 rings.

From all these ablation experiments, we decided to use the model with 5 rings and 90° as our *Base Model*. Considering both speed and accuracy, we decided to use the model with 3 rings and 90° as our *Light Model*.

Voting module vs. dilated convolution. Dilated convolution [54], which can include long-range features, could be considered as an alternative to our voting module. To compare performance, we trained models on `train2017` and evaluated them on `val2017` using the SS testing mode.

Baseline: We consider OAP with ResNet-101-DCN backbone as baseline. The last 1×1 convolution layer of center prediction branch in OAP, receives $H \times W \times D$ tensor and outputs object center heatmaps with a tensor of size $H \times W \times C$.

Baseline + Voting Module: We first adapt the last layer of center prediction branch in baseline to output $H \times W \times C \times R$ tensor, then attach our voting module on top of the center prediction branch. Adding the voting module increases parameters of the layer by R times. The log-polar vote field is 65×65 , and has 5 rings (90°). With 5 rings and 90° we end up with $R = 17$ regions.

Baseline + Dilated Convolution: We use dilated convolution with kernel size 4×4 and dilation rate 22 for the last layer of the center prediction branch in baseline. Using 4×4 kernel increases parameters 16 times which is approximately equal to R in the *Baseline + Voting Module*. Using dilation rate 22, the filter size becomes 67×67 which is close to 65×65 log-polar vote field.

For a fair comparison with *Baseline*, both *Baseline + Voting Module* and the *Baseline + Dilated Convolution* use ResNet-101-DCN backbone. Our voting module outperforms dilated convolution in all cases (Table 2).

4.2.2 Comparison with Baseline

In Table 3, we present the performance of HoughNet for different backbone networks, initializations and our base-vs-light model, on the `val2017` set. There is a significant speed difference between Base and Light models. Our light model with R-101-DCN backbone is the second fastest one (14.3 FPS) achieving 37.2 AP and 56.5 AP_{50} . We observe that initializing the backbone with a pretrained model improves the detection performance. In Table 4, we compare HoughNet’s performance with its baseline OAP [34] for two different backbones. HoughNet is especially effective for small objects, it improves the baseline by 2.1 and 2.2 AP points for R-101-DCN and HG-104 backbones, respectively. We also provide results for the recently introduced *moLRP* [19] metric, which combines localization, precision and recall in a single metric. Lower values are better.

4.2.3 Comparison with the State-of-the-art

For comparison with the state-of-the-art, we use Hourglass-104 [11] backbone. We train Hourglass model with a batch size of 36 for 100 epochs using the Adam optimizer [53]. We set the initial learning rate to 2.5×10^{-4} and divided it by 10 at epoch 90. Table 5 presents performances of HoughNet and several established state-of-the-art detectors. First, we compare HoughNet with OAP [34] since it is the model on which we built HoughNet. In OAP, they did not present any results for “from scratch” training. Instead they fine-tuned their model from ExtremeNet weights. When we do the same (i.e. initialize HoughNet with ExtremeNet weights), we obtain better results than OAP. However as expected, HoughNet is slower than OAP. Among the one-stage bottom-up object detectors, HoughNet performs on-par with the best bottom-up object detector by achieving 46.4 AP against 47.0 AP of CenterNet [16]. HoughNet outperforms CenterNet on AP_{50} (65.1 AP_{50} vs. 64.5 AP_{50}). Note that, since our model is initialized with ExtremeNet weights, which makes use of the segmentation masks in its own training, our model effectively uses more data compared to CenterNet. HoughNet is the fastest among one-stage bottom-up detectors. It is faster than CenterNet, CornerNet and more than twice as fast as ExtremeNet.

TABLE 5: Comparison with the state-of-the-art on COCO *test-dev*. The methods are divided into three groups: two-stage, one-stage top-down and one-stage bottom-up. The best results are boldfaced separately for each group. Backbone names are shortened: R is ResNet, X is ResNeXt, F is FPN and HG is HourGlass.* indicates that the FPS values were obtained on the same AWS machine with a V100 GPU using the official repos in SS setup. The rest of the FPS are from their corresponding papers. F. R-CNN is Faster R-CNN.

Method	Backbone	Initialize	Train size	Test size	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	FPS
<i>Two-stage detectors:</i>											
R-FCN [55]	R-101	ImageNet	800×800	600×600	29.9	51.9	-	10.8	32.8	45.0	5.9
CoupleNet [56]	R-101	ImageNet	ori.	ori.	34.4	54.8	37.2	13.4	38.1	50.8	-
F. R-CNN+++ [3]	R-101	ImageNet	1000×600	1000×600	34.9	55.7	37.4	15.6	38.7	50.9	-
F. R-CNN [57]	R-101-F	ImageNet	1000×600	1000×600	36.2	59.1	39.0	18.2	39.0	48.2	5.0
Mask R-CNN [17]	X-101-F	ImageNet	1300×800	1300×800	39.8	62.3	43.4	22.1	43.2	51.2	11.0
Cascade R-CNN [58]	R-101	ImageNet	-	-	42.8	62.1	46.3	23.7	45.5	55.2	12.0
PANet [59]	X-101	ImageNet	1400×840	1400×840	47.4	67.2	51.8	30.1	51.7	60.0	-
<i>One-stage detectors:</i>											
<i>Top Down:</i>											
SSD [7]	VGG-16	ImageNet	512×512	512×512	28.8	48.5	30.3	10.9	31.8	43.5	-
YOLOv3 [29]	Darknet	ImageNet	608×608	608×608	33.0	57.9	34.4	18.3	35.4	41.9	20.0
DSSD513 [30]	R-101	ImageNet	513×513	513×513	33.2	53.3	35.2	13.0	35.4	51.1	-
RefineDet (SS) [60]	R-101	ImageNet	512×512	512×512	36.4	57.5	39.5	16.6	39.9	51.4	-
RetinaNet [5]	X-101-F	ImageNet	1300×800	1300×800	40.8	61.1	44.1	24.1	44.2	51.2	5.4
RefineDet (MS) [60]	R-101	ImageNet	512×512	≤2.25×	41.8	62.9	45.7	25.6	45.1	54.1	-
OAP (SS) [34]	HG-104	ExtremeNet	512×512	ori.	42.1	61.1	45.9	24.1	45.5	52.8	9.6*
FSAF (SS) [61]	X-101	ImageNet	1300×800	1300×800	42.9	63.8	46.3	26.6	46.2	52.7	2.7
FSAF (MS) [61]	X-101	ImageNet	1300×800	~≤2.0×	44.6	65.2	48.6	29.7	47.1	54.6	-
FCOS [33]	X-101-F	ImageNet	1300×800	1300×800	44.7	64.1	48.4	27.6	47.5	55.6	7.0*
FreeAnchor (SS) [62]	X-101-F	ImageNet	1300×960	1300×960	44.9	64.3	48.5	26.8	48.3	55.9	-
OAP (MS) [34]	HG-104	ExtremeNet	512×512	≤1.5×	45.1	63.9	49.3	26.6	47.1	57.7	-
FreeAnchor (MS) [62]	X-101-F	ImageNet	1300×960	~≤2.0×	47.3	66.3	51.5	30.6	50.4	59.0	-
<i>Bottom Up:</i>											
ExtremeNet (SS) [12]	HG-104	-	511×511	ori.	40.2	55.5	43.2	20.4	43.2	53.1	3.0*
CornerNet (SS) [11]	HG-104	-	511×511	ori.	40.5	56.5	43.1	19.4	42.7	53.9	5.2*
CornerNet (MS) [11]	HG-104	-	511×511	≤1.5×	42.1	57.8	45.3	20.8	44.8	56.7	-
ExtremeNet (MS) [12]	HG-104	-	511×511	≤1.5×	43.7	60.5	47.0	24.1	46.9	57.6	-
CenterNet (SS) [16]	HG-104	-	511×511	ori.	44.9	62.4	48.1	25.6	47.4	57.4	4.8*
CenterNet (MS) [16]	HG-104	-	511×511	≤1.8×	47.0	64.5	50.7	28.9	49.9	58.9	-
HoughNet (SS)	HG-104	-	512×512	ori.	40.8	59.1	44.2	22.9	44.4	51.1	6.4*
HoughNet (MS)	HG-104	-	512×512	≤1.8×	44.0	62.4	47.7	26.4	45.4	55.2	-
HoughNet (SS)	HG-104	ExtremeNet	512×512	ori.	43.1	62.2	46.8	24.6	47.0	54.4	6.4*
HoughNet (MS)	HG-104	ExtremeNet	512×512	≤1.8×	46.4	65.1	50.7	29.1	48.5	58.1	-

We provide visualization of votes for sample detections of HoughNet for qualitative visual inspection (Fig. 7). These detections clearly show that HoughNet is able to make use of long-range visual evidence.

4.2.4 Analysis

Here we analyse the effect of Hough voting from two aspects. First, we inspect the error sources of both HoughNet and its baseline, and compare the two. This lets us understand how the voting module has improved the baseline. Secondly, we collect votes cast on a large number of images and analyse them to deduce interactions among object classes

Error Sources. Here we use the recently introduced “localisation recall precision” (LRP) metric [19], [20], which provides us with componentwise errors. These errors and the relative improvements yielded by the voting module are shown in Table 6. The largest relative improvement (+4.4%) is obtained in the “false negative” (FN) component, which shows that the voting module increases the recall rate, that is, object instances normally missed by the baseline are successfully recovered with the integration of

TABLE 6: Effect of voting module on three major error types, namely, localisation (Loc), false positive (FP) and false negative (FN), as measured by the LRP metric [19], [20]. Relative improvement is calculated as $(B - V)/B$ where B and V are the LRP errors of baseline and “baseline+voting module”, respectively. Lower is better. Largest relative improvement is achieved in the FN component, which indicates that voting module increases the recall over the baseline.

Model	Loc	FP	FN
Baseline w R-101-DCN	17.1	27.1	50.2
+ Voting Module	17.1	26.2	48.0
Relative Improvement	0%	+3.3%	+4.4%

the voting module. This shows the importance of voting. “False positive” (FP) component is also improved (+3.3%), which indicates that object instances marked as background by the baseline are corrected by the voting module. These improvements come with no degradation in localisation performance.

Interaction among object classes. We build a $C \times C$ (C

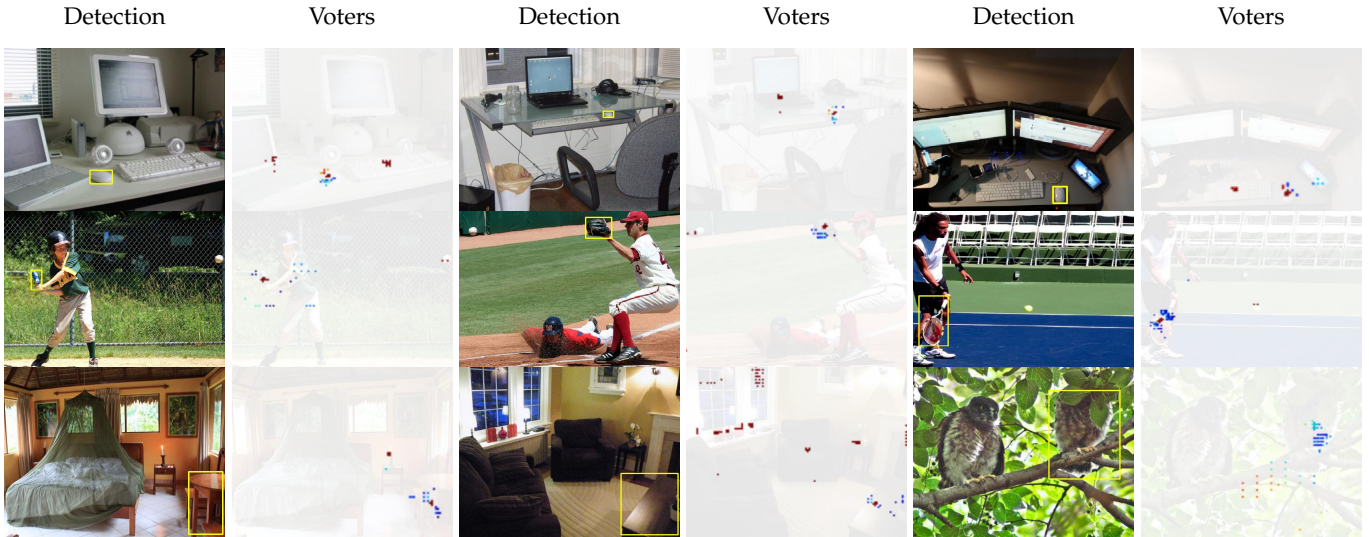


Fig. 7: Sample detections of HoughNet and their vote maps. In the “detection” columns, we show a correctly detected object, marked with a yellow bounding box. In the “voters” columns, the locations that vote for the detection are shown. Colors indicate vote strength based on the standard “jet” colormap (red is high, blue is low; Fig. 1). In the **top row**, there are three “mouse” detections. In all cases, in addition to the local votes (that are on the mouse itself), there are strong votes coming from nearby “keyboard” objects. This voting pattern is justified given that mouse and keyboard objects frequently co-appear. A similar behavior is observed in the detections of “baseball bat”, “baseball glove” and “tennis racket” in the **second row**, where they get strong votes from “ball” objects that are far-away. In the first example of the **bottom row**, “dining table” detection gets strong votes from the candle object, probably because they co-occur frequently. Candle is not among the 80 classes of COCO dataset. Similarly, in the second example in the **bottom row**, “dining table” has strong votes from objects and parts of a standard living room. In the last example, partially occluded bird gets strong votes (stronger than the local votes on the bird itself) from the tree branch.

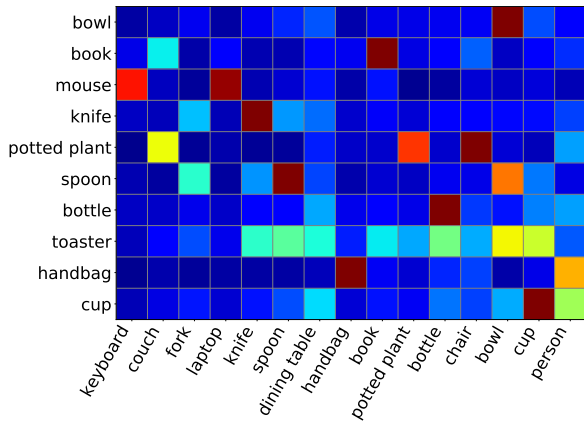


Fig. 8: Object detection voting activity among 10 vote-getter classes (rows) and 15 vote-giver classes (columns) on the COCO dataset. Color decodes voting activity (red: high, blue: low). Each detection in a vote-getter class might get votes from multiple object classes. For example, the *cup* objects (i.e. the “cup” row) get relatively high votes from *cup*, *person*, *dining table* and *bowl* classes. The full 80×80 matrix is provided in the supplementary material.

is the number of classes) matrix to visualize voting relations among classes on the COCO dataset using the R-101-DCN backbone of HoughNet. In this matrix, rows represent vote-getters and columns represent vote-givers. For each vote-getter class, we first identify center points of detections and find all locations that vote for these centers. Then, we sum

class probabilities of all the voter-giver locations and obtain a final C dimensional vector that summarizes the votes received by the vote-getter class. Fig. 8 shows a 10×15 matrix, which we selected from the full 80×80 matrix based on the following criterion: we selected the top 10 vote-getter classes with the maximum total votes. Next, for each vote-getter class, we selected the top 15 vote-giver classes that have the maximum voting activity. We provide the full 80×80 matrix in the supplementary material.

This matrix reflects many object co-occurrence relations and it captures object-to-object context well. For example, *mouse* class gets votes mostly from *mouse*, *keyboard* and *laptop* classes. *toaster* objects get votes from kitchen items. As a part of tableware items, both *spoon* and *knife* get votes from each other and other tableware items like *fork* and *bowl*. *Potted plant* objects get from indoor objects such as *chair* and *couch*.

4.2.5 A scalable approach to voting (independent of number of classes)

In HoughNet, there is a separate visual evidence tensor per object class (Fig. 2) and a voting process is run for each of these (Section 3.2). This linear dependence on the number of object classes might be problematic when the number of classes increases dramatically, which is the case for newer datasets such as the 1000-class LVIS dataset [63]. With this in mind, we designed a scalable variant of our voting module where, instead of having a separate visual evidence tensor per class, we create N tensors with $N \ll C$ (Fig. 9). In this design, these N tensors contain “visual evidence” scores that are shared among and common to all classes. As a

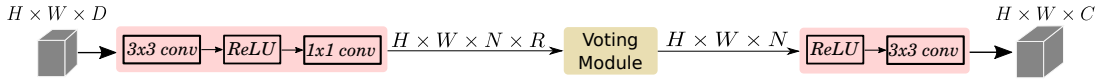


Fig. 9: A scalable variant of our voting process. Instead of having a separate visual evidence tensor for each of the C classes, we create N ($N \ll C$) visual evidence tensors that are shared and common to all classes.

TABLE 7: Ablation experiments for the scalable voting module. Models are trained on COCO `minitrain` and results are presented on SS testing mode. Using 8 visual evidence tensors yields the best result for all AP metrics.

N	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
4	24.9	40.7	25.9	7.7	27.5	37.5
8	25.9	42.3	26.7	8.5	28.4	39.1
12	24.6	40.6	25.6	7.8	27.0	36.6
16	25.0	41.1	25.8	8.2	27.4	37.9

result, the voting process is carried out N times instead of C times. We convert the resulting $H \times W \times N$ dimensional voting maps to $H \times W \times C$ dimensional object presence maps using a ReLU layer followed by a convolutional layer with 3×3 filters. We illustrate the scalable voting module in Fig. 9.

To analyze the scalable approach, we conducted ablation experiments for different values of N . We used the vote field from the *Light Model* setup with 3 rings and 90° bins, using the ResNet-101-DCN backbone. Models are trained on COCO `minitrain` and evaluated on `val2017` set in SS testing mode. Results are presented in Table 7. The model with $N = 8$ performs best among others.

To compare the scalable model with the light model, we trained it on COCO `train2017` and obtained results on `val2017`. The scalable model is almost 2 times faster than the light model (26.3 FPS vs. 14.3 FPS) in SS testing mode. However, it performs 1 AP point worse than the base and light models with the same backbone, obtaining 36.3 AP .

4.3 Hough Voting for Other Visual Detection Tasks

4.3.1 Video Object Detection

We conducted our experiments on the ImageNet VID dataset [64]. ImageNet VID dataset has 30 object categories. For training and evaluation, we followed widely adopted protocols [46], [49] and trained our models on 3,862 videos in the training set and evaluated their performances on 555 videos in the validation set using mean average precision (mAP) as the evaluation metric. In addition to the overall mAP, we also present mAP results for slow, medium, and fast groups of videos as done in previous work [46]. During the training of temporal models, we use a pair of video frames; a reference frame at time t and a nearby auxiliary frame at time $t + \delta$ where δ is randomly picked from the set $\{-4, \dots, +4\}$. At inference, we use fixed δ values and present our results for $\delta = \{-4\}$ and $\delta = \{-4, +4\}$ settings.

First, we obtained single-frame baseline results for our baseline (OAP [34]) and HoughNet. We trained OAP with 140 epochs using their official repository. For faster analysis, we trained our models with 80 epochs and divided the initial learning rate 1.25×10^{-4} by 10 at epoch 50. All models

TABLE 8: Results of video object detection on ImageNet VID validation set. Results are obtained without any test time augmentation. δ corresponds to the time offset of the auxiliary frame during inference.

Method	δ	mAP	mAP_F	mAP_M	mAP_S
<i>Single-frame models:</i>					
Baseline (OAP)	-	65.0	41.4	61.9	76.4
+ Voting (HoughNet)	-	68.8	45.8	66.1	79.1
<i>Temporal voting models</i>					
+ Feat. Agg.	-4	71.5	46.6	67.7	82.1
+ Motion Features	-4	73.9	50.4	71.5	82.8
+ Motion Features	-4,+4	74.9	51.5	72.5	83.6

are trained with a batch size of 32. Results in Table 8 show that HoughNet outperforms its baseline (68.8 vs 65.0) in this setting. Later, we conducted an ablation experiment to compare the performance of proposed temporal voting with feature aggregation over reference and auxiliary frames. To this end, both reference and auxiliary frames are passed through the backbone, then the output feature maps are concatenated and fed to the visual evidence prediction branch. This model improved the single-frame baseline result of HoughNet by 2.7 mAP points. Next, we experimented with temporal voting using motion features. Temporal voting improved the performance further by 2.4 points. Using two auxiliary frames (δ is $\{-4, +4\}$), we obtained an even better result achieving 74.9 mAP.

4.3.2 Instance Segmentation

In order to show the effectiveness of voting for instance segmentation, we first extended our baseline OAP to perform instance segmentation. Inspired by the recent method BlendMask [65], we added to OAP a new *prototype mask prediction* branch to predict category independent, single prototype mask. This branch outputs a $H \times W \times 1$ -dimensional feature map. In order to predict instance specific masks, we added another *attention map prediction* branch which outputs $H \times W \times 196$ dimensional attention features. These new branches have the same layer structure as other branches (see Fig. 2).

During training, we first get the *instance specific local prototypes* for each instance by cropping them from the prototype mask according to their location and box size. Next, we extract 1×196 *instance specific attention maps* using the center points of instances, from the predicted attention map. Later we obtain 14×14 instance specific attention maps by reshaping the attention features. We apply sigmoid normalization for both instance specific local prototypes and attention maps, and finally blend them by applying element-wise product. We use binary cross entropy as our loss function for segmentation. To extend HoughNet for the instance segmentation task, we add the new branches as

TABLE 9: Effect of voting module for the instance segmentation task on COCO val_{2017} set. Results are shown for both COCO segmentation and box AP. Models are trained on COCO train_{2017} . Results are obtained without any test time augmentation.

Method	AP^{seg}	AP_{50}^{seg}	AP_{75}^{seg}	AP_S^{seg}	AP_M^{seg}	AP_L^{seg}	AP^{box}	AP_{50}^{box}	AP_{75}^{box}	AP_S^{box}	AP_M^{box}	AP_L^{box}
Baseline	27.2	46.4	28.0	8.6	31.3	43.9	33.9	51.3	36.5	14.7	39.3	50.0
+ Voting	28.4	48.0	28.8	9.1	32.1	46.1	35.0	52.9	37.6	15.0	40.0	52.2

TABLE 10: Effect of voting module for the 3D object detection task on KITTI. Mean AP values of 5 models are presented with their standard deviations. Results are obtained without any test time augmentation.

Method	AP_e	AP_m	AP_h	AOS_e	AOS_m	AOS_h	BEV AP_e	BEV AP_m	BEV AP_h
Baseline	90.2±1.2	80.4±1.4	71.1±1.6	85.3±1.7	75.0±1.6	66.2±1.8	31.4±3.7	26.5±1.6	23.8±2.9
+ Voting	89.2±0.4	80.6±3.0	70.0±0.2	86.0±0.8	77.0±3.1	66.5±0.6	35.0±0.8	28.0±0.8	26.1±0.7

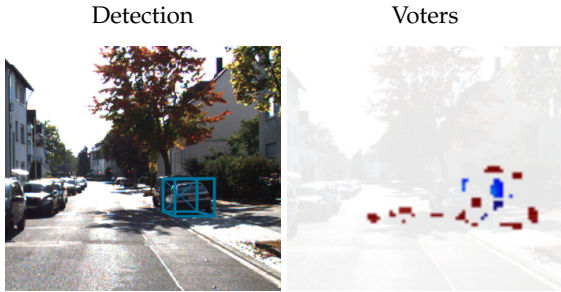


Fig. 10: Sample *car* detection of HoughNet from KITTI dataset and its vote map. We show a correctly detected object, marked with a blue 3D bounding box. In addition to the local votes, strong votes come from road and buildings.

described above and follow the same training process. We provide the network architecture of instance segmentation for both baseline and HoughNet in the supplementary material.

We present the instance segmentation results on the COCO dataset in Table 9. Both baseline and HoughNet models are trained for 80 epochs with a batch size of 32. Initial learning rate 1.25×10^{-4} is divided by 10 at epoch 50. Voting module of HoughNet improves instance segmentation performance by 1.2 AP points and outperforms the baseline for all instance segmentation and box AP metrics.

Integrating Hough voting to Mask R-CNN: Mask R-CNN [17] is an established state-of-the-art model for instance segmentation. It is an extension of the two-stage Faster R-CNN [6], where a new mask branch is added to predict object masks. Mask R-CNN [17] consists of a backbone and a Feature Pyramid Network (FPN) [57]. Each FPN level outputs bounding boxes, class and mask predictions separately. Attaching our voting module at the end of each FPN level would be costly in terms of inference time. Moreover, it may require searching for the optimal parameters of the vote field for each FPN level. To overcome these issues, we follow Libra R-CNN [66] which first resizes the multi-level FPN features to an intermediate size and then sums them to obtain a combined feature map. Libra R-CNN refines these combined features using a non-local neural network block [39]. In our experiments, we replace this non-local block with our voting module. Finally, refined features are rescaled and summed up with the output of each FPN

TABLE 11: Effect of voting module when integrated to Mask R-CNN for the instance segmentation task on COCO val_{2017} set. Results are shown for both COCO segmentation and box AP. Models are trained on COCO train_{2017} . Baseline results are obtained using the publicly available Mask R-CNN model from MMDetection [67].

Method	AP	AP_{50}	AP_{75}	AP^{seg}	AP_{50}^{seg}	AP_{75}^{seg}
Baseline	38.0	58.6	41.4	34.4	55.1	36.7
+ Voting	39.1	60.0	42.5	35.6	56.6	37.7

level to strengthen the original features. We conducted our experiments using ResNet-50 [3] with FPN backbone. We use MMDetection [67] framework with Pytorch [68] to implement our model. We present instance segmentation results in Table 11. Our voting module improves detection and instance segmentation performances of Mask R-CNN by 1.1 AP and 1.2 AP^{seg} points respectively, and outperforms the baseline for all AP metrics.

4.3.3 3D Object Detection

Here we conduct experiments in 3D object detection to see whether our voting module is useful. In addition to the class and location prediction as done in 2D object detection, a 3D object detector has to predict additional *depth*, *3D dimension* and *orientation* attributes. For this task also, we consider OAP [34] as our baseline. In addition to the branches from 2D object detection, OAP adds separate branches for these additional three attributes.

To show the effectiveness of the voting module, similar to 2D object detection, we attach our voting module with 3 rings and 90° to the class prediction branch. We experimented with *car* classes of KITTI dataset [69] using the training and validation splits from SubCNN [70]. Following OAP, we use the original image resolution 1280×384 both during training and inference. We trained the network with a batch size of 16 for 70 epochs with Adam optimizer [53]. Initial learning rate 1.25×10^{-4} was divided by 10 at epochs 45 and 60. For fair comparison with the baseline, we use Deep Layer Aggregation (DLA) [71] backbone as in OAP. More details related to both training and inference could be found in OAP [34].

Table 10 shows results for bounding box AP, average orientation score (AOS) and bird-eye-view bounding box

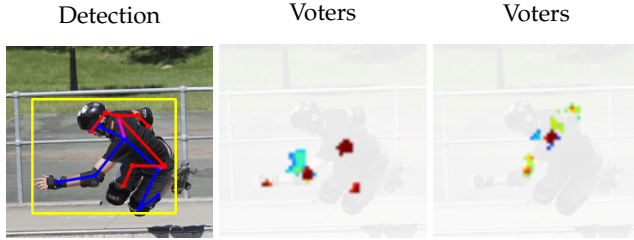


Fig. 11: Sample keypoint detection of HoughNet and their vote maps. In the “detection” column, we show a correctly detected object and its pose. Detection box is marked with a yellow bounding box, and pose estimation is shown with blue color for the left parts, red color for the right parts. Votes for the *left-elbow* and *left-shoulder* are shown respectively. The *left-elbow* detection gets strong votes from the *left-shoulder*, *left-knee* and *left-wrist*. The *left-shoulder* detection gets votes from the area spanning from *right-elbow* to *left-elbow*.

AP (BEV AP) at 3 levels of difficulty, namely, *easy*, *medium* and *hard*. AP values in the table correspond to the average AP at IoU threshold 0.5 at 11 recalls from 0.0 to 1.0 with a step size of 0.1. As also identified by OAP, since the recall threshold is small, the evaluation measures fluctuate up to 10% AP. To smooth the effect of fluctuations we trained 5 models and reported the averages together with standard deviation for each metric as in OAP.

The model with our voting module performs on par with the baseline for bounding box AP, and it outperforms baseline on AOS and BEV AP. Especially BEV AP is significantly improved and is more stable with much less variance compared to baseline. In Fig. 10, we show visualization of votes for sample *car* detections.

4.3.4 Keypoint detection: 2D Human Pose Estimation

In this section, we show the effectiveness of our voting module for 2D Human Pose Estimation task. Human Pose Estimation is the problem of detecting human joints (i.e. keypoints) in images. It is another detection task that our voting module could help by leveraging the interactions both among keypoints and other visual parts.

Our baseline (OAP [34]) considers pose estimation as a regression task and defines each keypoint with an offset to the center of the instance and directly regresses them in a separate branch. In order to refine keypoints, OAP also employs the common bottom-up multi-human pose estimation approach as in [72], [73], [74] and predicts heatmaps for each of k human joints in another branch. More detail on training and inference processes can be found in OAP [34].

We analyze the effect of voting by attaching our voting module to the class prediction and the keypoint heatmap prediction branches both separately and at the same time. We conducted our experiments on COCO dataset which has 17 keypoints for person object instances. For fair comparison with the baseline, we use DLA backbone and initialize it with the center detection model of OAP. We followed exactly the same training setup as in Section 4.2.1. The models are trained on `train2017` and evaluated on `val2017`. In Table 12, we show the results for both keypoint AP (AP^{kp})

and box AP (AP^{box}). Baseline results are obtained using the publicly available equivalent model trained with 140 epochs from the official repository of OAP. In all cases, attaching our voting module improved the baseline. Using the voting module in classification and keypoint heatmap prediction branches at the same time gives the best performance for both keypoint AP^{kp} and box AP^{box} . Attaching the voting module improves the baseline by 2.2 and 2.9 points for AP^{kp} and AP^{box} , respectively. Especially in AP^{box} , our voting module dramatically improved baseline by 7.8 points. We provide visualization of votes for sample detections in Fig. 11

4.4 Comparison with Non-local Neural Networks

A fair comparison between our voting module and non-local neural networks (NLNN) is not trivial. Our voting module is intended to work with a classification head to aggregate class conditional evidence whereas non-local neural blocks are targeted to be integrated into backbones and necks (e.g. FPN [57]) to operate at lower resolutions. Therefore, for a fair comparison, we integrated both models into two different locations: (1) a classification head and (2) a detection neck, that is the Feature Pyramid Network (FPN) layer.

Our first set of experiments aim to compare NLNN and our voting module when they are integrated to a classification head. To this end, we add one non-local block right before the last layer of OAP’s [34] classification head and compare its results with HoughNet. Since the output of the classification head has higher spatial resolution compared to backbone features, NLNN’s memory usage and inference time increase significantly. In our experiments, we use ResNet-101-DCN backbone and present results for single scale testing without any test time augmentations. As could be seen from Table 13, compared to NLNN, our voting module improves the baseline 0.4 AP points more, while being almost 2 times faster and consuming about 5 times less memory.

Next, we compare our voting module and non-local neural networks when both are integrated to a detection neck. Our comparison is conducted based on Libra R-CNN [66]. Libra R-CNN by default adopts a non-local neural block for its feature refinement step at the end of the FPN. We replaced the non-local neural block in the feature refinement layer with our voting module. Experiments are conducted on the ResNet-50 backbone with FPN. We used Pytorch [68] and MMDetection [67] for our implementation and followed the same training procedure as in Libra R-CNN. We present the results in Table 14. The baseline model is the base Libra R-CNN model with IoU-balanced sampling where balanced L1 loss is ignored (see [66] for the details). The model with our voting module performs on par with the Libra R-CNN model with a non-local block.

4.5 Hough Voting for an Image Generation Task

Another task where long-range interactions could be useful is the task of image generation from a given label map. There are two main approaches to solve this task; using unpaired and paired data for training. We take CycleGAN [75] and Pix2Pix [76] as our baselines for unpaired and paired

TABLE 12: Comparing our voting module with baseline for 2D human pose estimation on COCO val_{2017} set. The voting module is attached to the person classification and keypoint estimation branches separately and concurrently. Results are shown for both COCO keypoint and box AP. Models are trained on COCO $train_{2017}$. Results are presented on SS testing mode.

Method	AP^{kp}	AP_{50}^{kp}	AP_{75}^{kp}	AP_M^{kp}	AP_L^{kp}	AP^{box}	AP_{50}^{box}	AP_{75}^{box}	AP_S^{box}	AP_M^{box}	AP_L^{box}
Baseline	54.7	81.7	59.4	49.0	64.6	47.5	63.9	52.8	15.9	65.8	79.3
+ Voting for Person Class.	56.9	81.6	61.9	51.3	67.9	50.1	71.4	54.1	16.9	64.7	79.3
+ Voting for Keypoint Est.	56.8	81.5	61.2	50.7	68.0	50.2	70.9	54.3	17.2	64.4	79.4
+ Voting for Both	56.9	81.6	62.1	50.8	68.4	50.4	71.7	54.4	17.0	64.9	79.8

TABLE 13: Comparing our voting module and NLNN for object detection on OAP [34]. Models are trained on the COCO $train_{2017}$ set and results are presented on the COCO val_{2017} set.

Method	AP	AP_{50}	AP_{75}	FPS	Memory (GB)
OAP + NLNN	35.3	54.9	37.8	11	7.1
OAP + Voting	35.7	54.9	38.3	19	1.5

TABLE 14: Comparing our voting module and NLNN for object detection on Libra R-CNN. Models are trained on the COCO $train_{2017}$ set and results are presented on the COCO val_{2017} set.

Method	AP	AP_{50}	AP_{75}	FPS	Memory (GB)
Baseline	36.8	58.0	40.0	-	-
+ NLNN	37.7	59.4	40.9	11.3	2.0
+ Voting	37.8	58.5	41.2	11.5	1.8

approaches, respectively. We attach our voting module at the end of CycleGAN [75] and Pix2Pix [76] models.

For quantitative comparison, we use the Cityscapes [77] dataset. In Table 15, we present FCN scores [78] (which is used as the measure of success in this task) of CycleGAN and Pix2Pix with and without our voting module. To obtain the “without” result, we used the already trained model shared by the authors. We obtained the “with” result using the official training code from their repositories. In both cases evaluation was done using the official test and evaluation scripts from their repos. Results show that using the voting module improves FCN scores by large margins. Qualitative inspection also shows that when our voting module is attached, the generated images conform to the given input segmentation maps better (Fig. 12). This is the main reason for the quantitative improvement. Since Pix2Pix is trained with paired data, generated images follow input segmentation maps, however, Pix2Pix fails to generate small details.

TABLE 15: Comparison of FCN Scores for the “Labels to Photo” Task on the Cityscapes Dataset [77]. Higher scores are better.

Method	Per-pixel acc.	Per-class acc.	Class IOU
CycleGAN	0.43	0.14	0.09
+ Voting	0.52	0.17	0.13
pix2pix	0.71	0.25	0.18
+ Voting	0.76	0.25	0.20

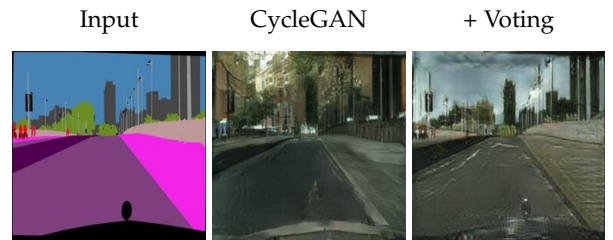


Fig. 12: Sample qualitative results for the “labels to photo” task. When integrated with CycleGAN, our voting module helps generate better images in the sense that the image conforms to the input label map better. CycleGAN fails to generate sky and buildings. For more visual results including Pix2Pix results please refer to [18] or supplementary material.

5 CONCLUSION

In this paper, we presented HoughNet, a new, one-stage, anchor-free, voting-based, bottom-up object detection method. HoughNet determines the presence of an object at a specific location by the sum of the votes cast on that location. Voting module of HoughNet is able to use both short and long-range evidence through its log-polar vote field. Thanks to this ability, HoughNet generalizes and enhances current object detection methodology, which typically relies on only local (short-range) evidence. We show that HoughNet performs on-par with the state-of-the-art bottom-up object detectors, and obtains comparable results with one-stage and two-stage methods. To further validate our proposal, we used the voting module of HoughNet in instance segmentation, 3D object detection, human pose estimation and “labels to photo” image generation tasks, and showed that our voting module consistently improves the baseline performances.

ACKNOWLEDGMENTS

This work was supported by the Scientific and Technological Research Council of Turkey (TÜBİTAK) through the project titled “Object Detection in Videos with Deep Neural Networks” (grant #117E054). The numerical calculations reported in this paper were partially performed at TÜBİTAK ULAKBİM, High Performance and Grid Computing Center (TRUBA resources). We also gratefully acknowledge the support of the AWS Cloud Credits for Research program.

REFERENCES

- [1] R. B. Girshick, P. F. Felzenszwalb, and D. McAllester, "Discriminatively trained deformable part models, release 5," <http://people.cs.uchicago.edu/~rbg/latent-release5/>.
- [2] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–1645, 2009.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [4] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *European Conf. on Computer Vision*. Springer, 2014, pp. 740–755.
- [5] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *IEEE International Conference on Computer Vision*, 2017.
- [6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 91–99.
- [7] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European Conf. on Computer Vision*, 2016.
- [8] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [9] B. Leibe, A. Leonardis, and B. Schiele, "Robust object detection with interleaved categorization and segmentation," *International Journal of Computer Vision*, vol. 77, no. 1, pp. 259–289, May 2008.
- [10] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part based models," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2010.
- [11] H. Law and J. Deng, "Cornersnet: Detecting objects as paired keypoints," in *European Conf. on Computer Vision*, 2018, pp. 734–750.
- [12] X. Zhou, J. Zhuo, and P. Krähenbühl, "Bottom-up object detection by grouping extreme and center points," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [13] P. V. C. Hough, "Machine Analysis of Bubble Chamber Pictures," vol. C590914, pp. 554–558, 1959.
- [14] D. H. Ballard *et al.*, "Generalizing the hough transform to detect arbitrary shapes," *Pattern Recognition*, 1981.
- [15] M. Land and B. Tatler, *Looking and acting: vision and eye movements in natural behaviour*. Oxford University Press, 2009.
- [16] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *IEEE International Conference on Computer Vision*, 2019.
- [17] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask r-cnn," *IEEE International Conference on Computer Vision*, pp. 2980–2988, 2017.
- [18] N. Samet, S. Hicsonmez, and E. Akbas, "Houghnet: Integrating near and long-range evidence for bottom-up object detection," in *ECCV*. Springer, 2020, pp. 406–423.
- [19] K. Oksuz, B. Cam, E. Akbas, and S. Kalkan, "Localization recall precision (LRP): A new performance metric for object detection," in *European Conference on Computer Vision (ECCV)*, 2018.
- [20] K. Oksuz, B. C. Cam, S. Kalkan, and E. Akbas, "One metric to measure them all: Localisation recall precision (LRP) for evaluating visual detection tasks," *under review at TPAMI*, 2020.
- [21] S. Belongie, J. Malik, and J. Puzicha, "Shape matching and object recognition using shape contexts," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 24, no. 4, pp. 509–522, 2002.
- [22] E. Akbas and M. P. Eckstein, "Object detection through search with a foveated visual system," *PLoS computational biology*, vol. 13, no. 10, p. e1005743, 2017.
- [23] V. J. Traver and A. Bernardino, "A review of log-polar imaging for visual perception in robotics," *Robotics and Autonomous Systems*, vol. 58, no. 4, pp. 378–398, 2010.
- [24] R. Okada, "Discriminative generalized hough transform for object detection," in *IEEE International Conference on Computer Vision*, 2009.
- [25] J. Gall and V. Lempitsky, "Class-specific hough forests for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- [26] N. Razavi, J. Gall, P. Kohli, and L. Van Gool, "Latent hough transform for object detection," in *European Conf. on Computer Vision*, 2012.
- [27] S. Maji and J. Malik, "Object detection using a max-margin hough transform," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- [28] O. Barinova, V. Lempitsky, and P. Kholi, "On detection of multiple object instances using hough transforms," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 34, no. 9, pp. 1773–1784, 2012.
- [29] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [30] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "Dssd: Deconvolutional single shot detector," *arXiv preprint arXiv:1701.06659*, 2017.
- [31] J. Ren, X. Chen, J. Liu, W. Sun, J. Pang, Q. Yan, Y.-W. Tai, and L. Xu, "Accurate single stage detector using recurrent rolling convolution," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [32] J. Wang, K. Chen, S. Yang, C. C. Loy, and D. Lin, "Region proposal by guided anchoring," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [33] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *IEEE International Conference on Computer Vision*, 2019.
- [34] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," in *arXiv preprint arXiv:1904.07850*, 2019.
- [35] C. R. Qi, O. Litany, K. He, and L. J. Guibas, "Deep hough voting for 3d object detection in point clouds," in *IEEE International Conference on Computer Vision*, 2019.
- [36] A. Sheshkus, A. Ingacheva, V. Arlazarov, and D. Nikolaev, "Houghnet: neural network architecture for vanishing points detection," 2019.
- [37] E. Gabriel, M. Schleiss, H. Schramm, and C. Meyer, "Analysis of the discriminative generalized hough transform as a proposal generator for a deep network in automatic pedestrian and car detection," *Journal of Electronic Imaging*, vol. 27, no. 5, p. 051228, 2018.
- [38] I. Lifshitz, E. Fetaya, and S. Ullman, "Human pose estimation using deep consensus voting," in *European Conf. on Computer Vision*, 2016.
- [39] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803.
- [40] H. Hu, J. Gu, Z. Zhang, J. Dai, and Y. Wei, "Relation networks for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3588–3597.
- [41] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," *European Conf. on Computer Vision*, 2020.
- [42] H. Qiu, Y. Ma, Z. Li, S. Liu, and J. Sun, "Borderdet: Border feature for dense object detection," in *European Conf. on Computer Vision*, 2020, pp. 549–564.
- [43] C. Chi, F. Wei, and H. Hu, "Relationnet++: Bridging visual representations for object detection via transformer decoder," *Advances in Neural Information Processing Systems*, 2020.
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [45] H. Zhao, J. Jia, and V. Koltun, "Exploring self-attention for image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 076–10 085.
- [46] X. Zhu, Y. Wang, J. Dai, L. Yuan, and Y. Wei, "Flow-guided feature aggregation for video object detection," in *ICCV*, 2017.
- [47] H. Deng, Y. Hua, T. Song, Z. Zhang, Z. Xue, R. Ma, N. Robertson, and H. Guan, "Object guided external memory network for video object detection," in *ICCV*, 2019.
- [48] C. Guo, B. Fan, J. Gu, Q. Zhang, S. Xiang, V. Prinet, and C. Pan, "Progressive sparse local attention for video object detection," in *ICCV*, 2019.
- [49] Y. Chen, Y. Cao, H. Hu, and L. Wang, "Memory enhanced global-local aggregation for video object detection," in *CVPR*, 2020.
- [50] G. Bertasius, C. Feichtenhofer, D. Tran, J. Shi, and L. Torresani, "Learning temporal pose estimation from sparsely labeled videos," in *NIPS*, 2019.
- [51] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis, "Soft-nms—improving object detection with one line of code," in *IEEE International Conference on Computer Vision*, 2017, pp. 5561–5569.

- [52] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," in *European Conf. on Computer Vision*, 2018, pp. 466–481.
- [53] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [54] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *CoRR*, 2015.
- [55] J. Dai, Y. Li, K. He, and J. Sun, "R-FCN: Object detection via region-based fully convolutional networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 379–387.
- [56] Y. Zhu, C. Zhao, J. Wang, X. Zhao, Y. Wu, and H. Lu, "Couplenet: Coupling global structure with local parts for object detection," in *IEEE International Conference on Computer Vision*, 2017, pp. 4126–4134.
- [57] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 936–944.
- [58] Z. Cai and N. Vasconcelos, "Cascade R-CNN: Delving into high quality object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6154–6162.
- [59] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8759–8768.
- [60] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Single-shot refinement neural network for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4203–4212.
- [61] C. Zhu, Y. He, and M. Savvides, "Feature selective anchor-free module for single-shot object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [62] X. Zhang, F. Wan, C. Liu, R. Ji, and Q. Ye, "Freeanchor: Learning to match anchors for visual object detection," in *Advances in Neural Information Processing Systems*, 2019.
- [63] A. Gupta, P. Dollar, and R. Girshick, "LVIS: A dataset for large vocabulary instance segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [64] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *IJCV*, 2015.
- [65] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, "Blend-mask: Top-down meets bottom-up for instance segmentation," in *CVPR*, 2020.
- [66] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra r-cnn: Towards balanced learning for object detection," in *CVPR*, 2019.
- [67] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin, "MMDetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.
- [68] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems*, 2019, pp. 8024–8035.
- [69] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 3354–3361.
- [70] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Subcategory-aware convolutional neural networks for object proposals and detection," in *IEEE Winter Conference on Applications of Computer Vision*, 2017, pp. 924–933.
- [71] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2403–2412.
- [72] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *Advances in Neural Information Processing Systems*, 2017, pp. 2277–2287.
- [73] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7291–7299.
- [74] G. Papandreou, T. Zhu, L.-C. Chen, S. Gidaris, J. Tompson, and K. Murphy, "Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model," in *European Conf. on Computer Vision*, 2018, pp. 269–286.
- [75] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *IEEE International Conference on Computer Vision*, 2017, pp. 2223–2232.
- [76] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1125–1134.
- [77] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3213–3223.
- [78] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.



Nermin Samet is a postdoctoral researcher at École des Ponts (ENPC) ParisTech, IMAGINE Research Group. She received her Ph.D. degree in 2021 from the Department of Computer Engineering, Middle East Technical University, Ankara, Turkey. Her M.Sc. and B.Sc. degrees in computer engineering are from Bilkent University and Hacettepe University respectively. Her primary research interest is computer vision focusing on object detection.



Samet Hicsonmez is a Ph.D. candidate at the Department of Computer Engineering, Hacettepe University, Ankara, Turkey. Before joining Hacettepe, he received his M.Sc. degree from TOBB Economics and Technology University, Computer Engineering Department in 2013, and B.Sc. degree from Istanbul Technical University, Electronics Engineering in 2009. His primary research interests are computer vision focusing on GANs and object detection.



Emre Akbas is an assistant professor at the Department of Computer Engineering, Middle East Technical University (METU), Ankara, Turkey. He received his Ph.D. degree from the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign in 2011. His M.Sc. and B.Sc. degrees in computer science are both from METU. Prior to joining METU, he was a postdoctoral research associate at the Department of Psychological and Brain Sciences, University of California Santa

Barbara. His research interests are in computer vision and deep learning with a focus on object detection and human pose estimation.