

Attend Refine Repeat: Active Box Proposal Generation via In-Out Localization

Spyros Gidaris
spyros.gidaris@enpc.fr

Nikos Komodakis
nikos.komodakis@enpc.fr

Université Paris-Est, École des Ponts
ParisTech
Paris, France

Abstract

The problem of computing category agnostic bounding box proposals is utilized as a core component in many computer vision tasks and thus has lately attracted a lot of attention. In this work we propose a new approach to tackle this problem that is based on an active strategy for generating box proposals that starts from a set of seed boxes, which are uniformly distributed on the image, and then progressively moves its attention on the promising image areas where it is more likely to discover well localized bounding box proposals. We call our approach *AttractionNet* and a core component of it is a CNN-based category agnostic object location refinement module that is capable of yielding accurate and robust bounding box predictions regardless of the object category.

We extensively evaluate our *AttractionNet* approach on several image datasets (i.e. COCO, PASCAL, ImageNet detection and NYU-Depth V2 datasets) reporting on all of them state-of-the-art results that surpass the previous work in the field by a significant margin and also providing strong empirical evidence that our approach is capable to generalize to unseen categories. Furthermore, we evaluate our *AttractionNet* proposals in the context of the object detection task using a VGG16-Net based detector and the achieved detection performance on COCO manages to significantly surpass all other VGG16-Net based detectors while even being competitive with a heavily tuned ResNet-101 based detector. Code as well as box proposals computed for several datasets are available at: <https://github.com/gidariss/AttractionNet>.

1 Introduction

Category agnostic object proposal generation is a computer vision task that has received an immense amount of attention over the last years. Its definition is that for a given image a small set of instance segmentations or bounding boxes must be generated that will cover with high recall all the objects that appear in the image regardless of their category. In object detection, applying the recognition models to such a reduced set of category independent location hypothesis [16] instead of an exhaustive scan of the entire image [10, 36], has the advantages of drastically reducing the amount of recognition model evaluations and thus allowing the use of more sophisticated machinery for that purpose. As a result, proposal based detection systems manage to achieve state-of-the-art results and have become the dominant paradigm in the object detection literature [2, 13, 14, 15, 16, 19, 54, 58, 42]. Object proposals have also been used in various other tasks, such as weakly-supervised object detection [2],

exemplar 2D-3D detection [60], visual semantic role labelling [17], caption generation [24] or visual question answering [82].

In this work we focus on the problem of generating bounding box object proposals rather than instance segmentations. Several approaches have been proposed in the literature for this task [1, 2, 6, 6, 8, 18, 25, 26, 29, 40, 42]. Among them our work is most related to the CNN-based objectness scoring approaches [12, 27, 53] that recently have demonstrated state-of-the-art results [52, 53].

In the objectness scoring paradigm, a large set of image boxes is ranked according to how likely it is for each image box to tightly enclose an object — regardless of its category — and then this set is post-processed with a non-maximum-suppression step and truncated to yield the final set of object proposals. In this context, Kuo *et al.* [27] with their Deep-Box system demonstrated that training a convolutional neural network to perform the task of objectness scoring can yield superior performance over previous methods that were based on low level cues and they provided empirical evidence that it can generalize to unseen categories. In order to avoid evaluating the computationally expensive CNN-based objectness scoring model on hundreds of thousands image boxes, which is necessary for achieving good localization of all the objects in the image, they use it only to re-rank the proposals generated from a faster but less accurate proposal generator thus being limited by its localization performance. Instead, more recent CNN-based approaches apply their models only to ten of thousands image boxes, uniformly distributed in the image, and jointly with objectness prediction they also infer the bounding box of the closest object to each input image box. Specifically, the Region Proposal Network in Faster-RCNN [34] performs bounding box regression for that purpose while the DeepMask method predicts the foreground mask of the object centred in the image box and then it infers the location of the object’s bounding box by extracting the box that tightly encloses the foreground pixels. The latter has demonstrated state-of-the-art results and was recently extended with a top-down foreground mask refinement mechanism that exploits the convolutional feature maps at multiple depths of a neural network [52].

Our work is also based on the paradigm of having a CNN model that given an image box it jointly predicts its objectness and a new bounding box that is better aligned on the object that it contains. However, we opt to advance the previous state-of-the-art in box proposal generation in two ways: (1) improving the object’s bounding box prediction step (2) actively generating the set of image boxes that will be processed by the CNN model.

Regarding the bounding box inference step we exploit the recent advances in object detection where Gidaris and Komodakis [14] showed how to improve the object-specific localization accuracy. Specifically, they replaced the bounding box regression step with a localization module, called *LocNet*, that given a search region it infers the bounding box of the object inside the search region by assigning membership probabilities to each row and each column of that region and they empirically proved that this localization task is easier to be learned from a convolutional neural network thus yielding more accurate box predictions during test time. Given the importance of having accurate bounding box locations in the proposal generation task, we believe that it would be of great interest to develop and study a *category agnostic* version of *LocNet* for this task.

Our second idea for improving the box proposal generation task stems from the following observation. Recent state-of-the-art box proposal methods evaluate only a relatively small set of image boxes (in the order of $10k$) uniformly distributed in the image and rely on the bounding box prediction step to fix the localization errors. However, depending on how far an object is from the closest evaluated image box, both the objectness scoring and the

bounding box prediction for that object could be imperfect. For instance, Hosang *et al.* [24] showed that in the case of the detection task the correct recognition of an object from an image box is correlated with how well the box encloses the object. Given how similar are the tasks of category-specific object detection and category-agnostic proposal generation, it is safe to assume that a similar behaviour will probably hold for the latter one as well. Hence, in our work we opt for an *active* object localization scheme, which we call *Attend Refine Repeat* algorithm, that starting from a set of seed boxes it progressively generates newer boxes that are expected with higher probability to be on the neighbourhood or to tightly enclose the objects of the image. Thanks to this localization scheme, our box proposal system is capable to both correct initially imperfect bounding box predictions and to give higher objectness score to candidate boxes that are more well localized on the objects of the image.

To summarize, our contributions with respect to the box proposal generation task are:

- We developed a box proposal system that is based on an improved category-agnostic object location refinement module and on an active box proposal generation strategy that behaves as an attention mechanism that focus on the promising image areas in order to propose objects. We call the developed box proposal system *AttractionNet: (Attend) (Refine) Repeat: (Active) Box Proposal Generation via (In-) (Out) Localization (Net)work*.
- We exhaustively evaluate our system both on PASCAL and on the more challenging COCO datasets and we demonstrate significant improvement with respect to the state-of-the-art on box proposal generation. Furthermore, we provide strong evidence that our object location refinement module is capable of generalizing to unseen categories by reporting results for the unseen categories of ImageNet detection task and NYU-Depth dataset.
- Finally, we evaluate our box proposal generation approach in the context of the object detection task using a VGG16-Net based detection system and the achieved average precision performance on the COCO test-dev set manages to significantly surpass all other VGG16-Net based detection systems while even being on par with the ResNet-101 based detection system of He *et al.* [21].

The remainder of the paper is structured as follows: We describe our box proposal methodology in section §2, we show experimental results in section §3 and we present our conclusions in section §4.

2 Our approach

2.1 Active bounding box proposal generation

The active box proposal generation strategy that we employ in our work, which we call *Attend Refine Repeat* algorithm, starts from a set of seed boxes, which only depend on the image size, and it then sequentially produces newer boxes that will better cover the objects of the image while avoiding the "objectless" image areas (see Figure 1). At the core of this algorithm lies a CNN-based box proposal model that, given an image I and the coordinates of a box B , executes the following operations:

Category agnostic object location refinement: this operation returns the coordinates of a new box \tilde{B} that would be more tightly aligned on the object near B . In case there are more than one objects in the neighbourhood of B then the new box \tilde{B} should be targeting the object closest to the input box B , where by closest we mean the object

Algorithm: Attend Refine Repeat

```

Input : Image  $\mathbf{I}$ 
Output: Bounding box proposals  $\mathbf{P}$ 
 $\mathbf{C} \leftarrow \emptyset, \mathbf{B}^0 \leftarrow$  seed boxes
for  $t \leftarrow 1$  to  $T$  do
  /* Attend & Refine procedure */
   $\mathbf{O}^t \leftarrow$  ObjectnessScoring( $\mathbf{B}^{t-1}|\mathbf{I}$ )
   $\mathbf{B}^t \leftarrow$  ObjectLocationRefinement( $\mathbf{B}^{t-1}|\mathbf{I}$ )
   $\mathbf{C} \leftarrow \mathbf{C} \cup \{\mathbf{B}^t, \mathbf{O}^t\}$ 
end
 $\mathbf{P} \leftarrow$  NonMaxSuppression( $\mathbf{C}$ )

```

that its bounding box has the highest intersection over union (IoU) overlap with the input box B .

Category agnostic objectness scoring: this operation scores the box B based on how likely it is to tightly enclose an object, regardless of its category.

The pseudo-code of the *Attend Refine Repeat* algorithm is provided in Algorithm 1. Specifically, it starts by initializing the set of candidate boxes \mathbf{C} to the empty set and then creates a set of seed boxes \mathbf{B}^0 by uniformly distributing boxes of various fixed sizes in the image (similar to Cracking Bing [43]). Then on each iteration t it estimates the objectness \mathbf{O}^t of the boxes generated in the previous iteration, \mathbf{B}^{t-1} , and it refines their location (resulting in boxes \mathbf{B}^t) by attempting to predict the bounding boxes of the objects that are closest to them. The results $\{\mathbf{B}^t, \mathbf{O}^t\}$ of those operations are added to the candidates set \mathbf{C} and the algorithm continues. In the end, non-maximum-suppression [44] is applied to the candidate box proposals \mathbf{C} and the top K box proposals, set \mathbf{P} , are returned.

The advantages of having an algorithm that sequentially generates new box locations given the predictions of the previous stage are two-fold:

- **Attention mechanism:** First, it behaves as an attention mechanism that, on each iteration, focuses more and more on the promising locations (in terms of box coordinates) of the image (see Figure 1). As a result of this, boxes that tightly enclose the image objects are more likely to be generated and to be scored with high objectness confidence.
- **Robustness to initial boxes:** Furthermore, it allows to refine some initially imperfect box predictions or to localize objects that might be far (in terms of center location, scale and/or aspect ratio) from any seed box in the image. This is illustrated via a few characteristic examples in Figure 2. As shown in each of these examples, starting from a seed box, the iterative bounding box predictions gradually converge to the closest (in terms of center location, scale and/or aspect ratio) object without actually being affected from any nearby instances.

2.2 CNN-based box proposal model

In this section we describe in more detail the object localization and objectness scoring modules of our box proposal model as well as the CNN architecture that implements the entire *Attend Refine Repeat* algorithm that was presented above.

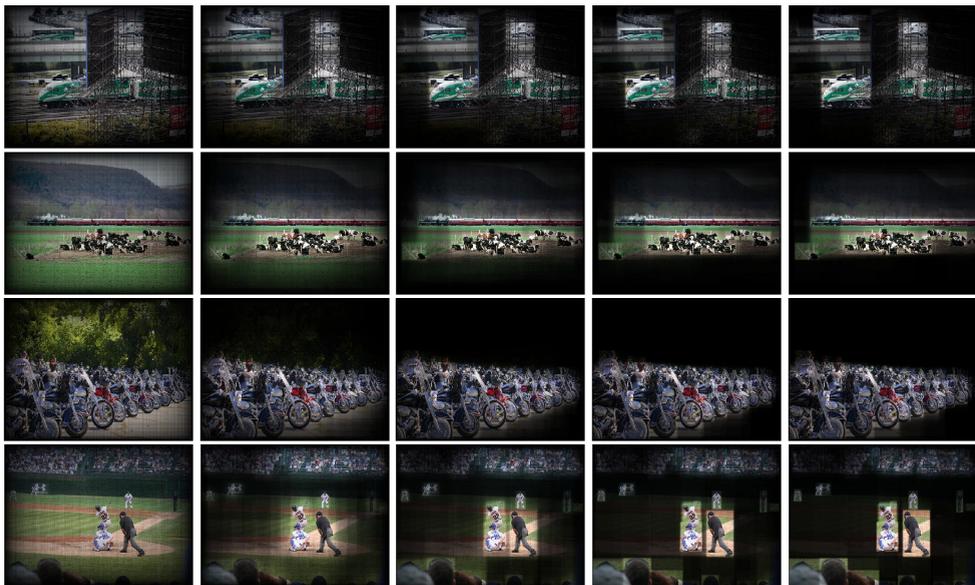


Figure 1: Illustration of the image areas being attended by our box proposal generator algorithm at each iteration. In the first iteration the box proposal generator attends the entire image since the seed boxes are created by uniformly distributing boxes across the image. However, as the algorithm progresses its attention is concentrated on the image areas that actually contain objects.

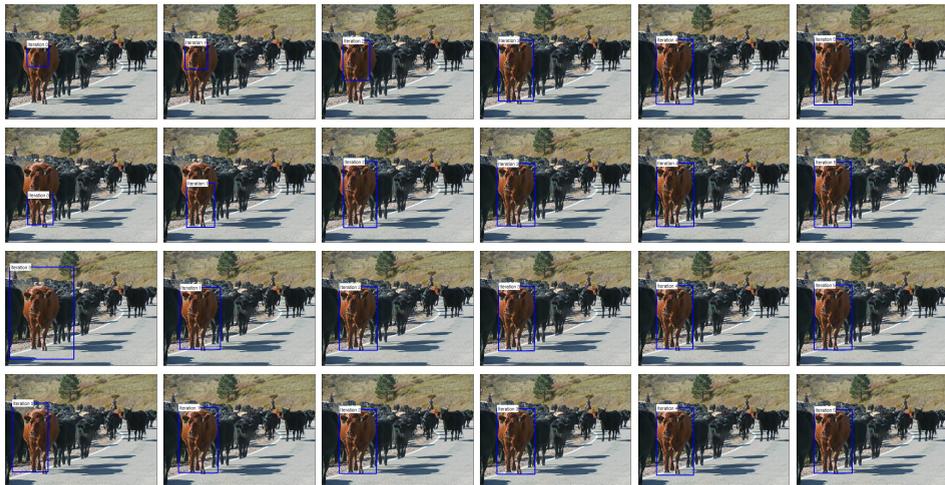


Figure 2: Illustration of the consecutive bounding box predictions made by our category agnostic location refinement module. In each row, from left to right we depict a seed box (iteration 0) and the bounding box predictions in each iteration. Despite the fact that the seed box might be quite far from the object (in terms of center location, scale and/or aspect ratio) the refinement module has no problem in converging to the bounding box closest to the seed box object. This capability is not affected even in the case that the seed box contains also other instances of the same category as in rows 3 and 4.

2.2.1 Object location refinement module

In order for our active box proposals generation strategy to be effective, it is very important to have an accurate and robust category agnostic object location refinement module. Hence we follow the paradigm of the recently introduced LocNet model [14] that has demonstrated superior performance in the category specific object detection task over the typical bounding box regression paradigm [13, 15, 16, 17] by formulating the problem of bounding box prediction as a dense classification task. Here we use a properly adapted version of that model for the task at hand.

At a high level, given as input a bounding box B , the location refinement module first defines a search region $R = \gamma B$ (i.e., the region of B enlarged by a factor γ) over which it is going to next search for a new refined bounding box. To achieve this, it considers a discretization of the search region R into M columns as well as M rows, and yields two probability vectors, $p_x = \{p_{x,i}\}_{i=1}^M$ and $p_y = \{p_{y,i}\}_{i=1}^M$, for the M columns and the M rows respectively of R , where these probabilities represent the likelihood of those elements (rows or columns) to be inside the target box B^* (these are also called *in-out* probabilities in the original LocNet model). Each time the target box B^* is defined to be the bounding box of the object closest to the input box B . Finally, given those *in-out* probabilities, the object location \tilde{B} inference is formulated as a simple maximum likelihood estimation problem that maximizes the likelihood of the *in-out* elements of \tilde{B} . A visual illustration of the above process through a few examples is provided in Fig. 3 (for further details about the LocNet model we refer the interested reader to [14]).

We note that in contrast to the original LocNet model that is optimized to yield a different set of probability vectors of each category in the training set, here our category-agnostic version is designed to yield a single set of probability vectors that should accurately localize any object regardless of its category (see also section §2.2.3 that describes in detail the overall architecture of our proposed model). It should be also mentioned that this is a more challenging task to learn since, in this case, the model should be able to localize the target objects even if they are in crowded scenes with other objects of the same appearance and/or texture (see the two left-most examples of Figure 3) without exploiting any category supervision during training that would help it to better capture the appearance characteristics of each object category. On top of that, our model should be able to localize objects of unseen categories. In the right-most example of Figure 3, we provide an indicative result produced by our model that verifies this test case. In this particular example, we apply a category-agnostic refinement module trained on PASCAL to an object whose category ("clock") was not present in the training set and yet our trained model had no problem of confidently predicting the correct location of the object. In section 3.2 of the paper we also provide quantitative results about the generalization capabilities of the location refinement module.

2.2.2 Objectness scoring module

The functionality of the objectness scoring module is that it gets as input a box B and yields a single probability p_{obj} of whether or not this box tightly encloses an object, regardless of what the category of that object might be.

2.2.3 AttractionNet architecture

We call the overall network architecture that implements the *Attend Refine Repeat* algorithm with its *In-Out* object location refinement module and its objectness scoring module, *At-*

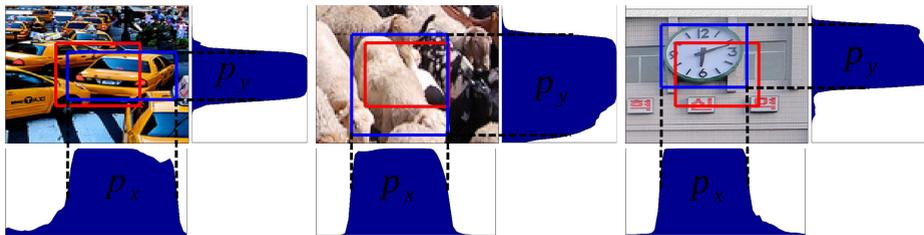


Figure 3: Illustration of the bounding box prediction process that is performed by our location refinement module. In each case the red rectangle is the input box B , the blue rectangle is the predicted box and the depicted image crops are the search regions where the refinement module "looks" in order to localize the target object. On the bottom and on the right side of the image crop we visualize the p_x and p_y probability vectors respectively that our location refinement module yields in order to localize the target object. Ideally, those probabilities should be equal to 1 for the elements (columns/rows) that overlap with the target box and 0 everywhere else. In the two left-most examples, the object location refinement module manages to correctly return the location of the target objects despite the fact that they are in crowded scenes with other objects of the same appearance and/or texture. In the right-most example, we qualitatively test the generalization capability of an location refinement module trained on PASCAL when applied to an object whose category ("clock") was not present in the training set.

*tractionNet*¹. Given an image I , our *AttractionNet* model will be required to process multiple image boxes of various sizes, by two different modules and repeat those processing steps for several iterations of the *Attend Refine Repeat* algorithm. So, in order to have an efficient implementation we follow the SPP-Net [14] and Fast-RCNN [15] paradigm and share the operations of the first convolutional layers between all the boxes, as well as across the two modules and all the *Attend Refine Repeat* algorithm repetitions (see Figure 4). Specifically, our *AttractionNet* model first forwards the image I through a first sequence of convolutional layers (conv. layers of VGG16-Net [39]) in order to extract convolutional feature maps F_l from the entire image. Then, on each iteration t the box-wise part of the architecture, which we call *Attend & Refine Network*, gets as input the image convolutional feature maps F_l and a set of box locations \mathbf{B}^{t-1} and yields the refined bounding box locations \mathbf{B}^t and their objectness scores \mathbf{O}^t using its object location refinement module sub-network and its objectness scoring module sub-network respectively. In Figure 5 we provide the work-flow of the *Attend & Refine Network* when processing a single input box B . The architecture of its two sub-networks is described in more detail in the rest of this section:

Object location refinement module sub-network. This module gets as input the feature map F_l and the search region R and yields the probability vectors p_x and p_y of that search region with a network architecture similar to that of LocNet. Key elements of this architecture is that it branches into two heads, the X and Y, each responsible for yielding the p_x or the p_y outputs. Differently from the original LocNet architecture, the convolutional layers of this sub-network output 128 feature channels instead of 512, which speeds up the processing by a factor of 4 without affecting the category-agnostic localization accuracy. Also, in order to yield a fixed size feature for the R region, instead of region adaptive max-pooling this sub-network uses region bilinear pooling [9, 23] that in our initial experiments gave slightly better results. Finally, our version is designed to yield two probability vectors of size M^2 , instead of $C \times 2$ vectors of size M (where C is the number of categories), given that in our

¹*AttractionNet* : (Att)end (R)efine Repeat: (Act)ive Box Proposal Generation via (I)n-(O)ut Localization (Net)work

²Here we use $M = 56$.

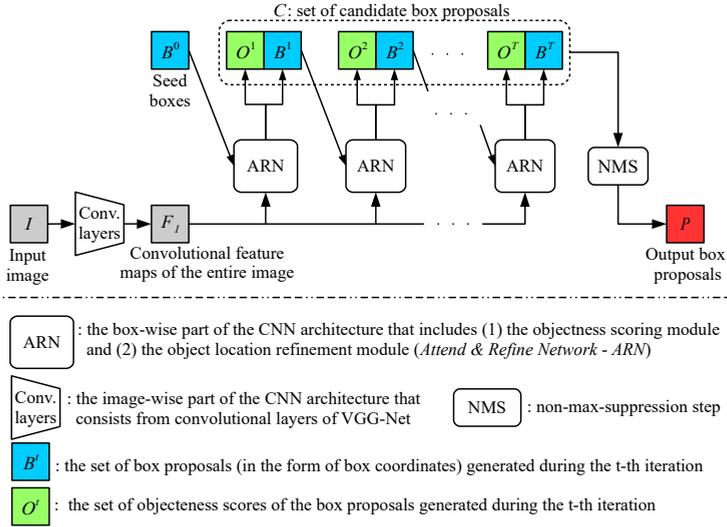


Figure 4: **AttracNet work-flow.** The *Attend Refine Repeat* algorithm is implemented through a CNN model, called *AttracNet*, whose run-time work-flow (when un-rolled over time) is illustrated here. On each iteration t the box-wise part of the architecture (*Attend & Refine Network: ARN*) gets as input the image convolutional feature maps F_t (extracted from the image-wise part of the CNN architecture) as well as a set of box locations B^{t-1} and yields the refined bounding box locations B^t and their objectness scores O^t using its *category agnostic object location refinement* module and its *category agnostic objectness scoring* module respectively. To avoid any confusion, note that our *AttracNet* model does not include any recurrent connections.

case we aim for category-agnostic object location refinement.

Objectness scoring module sub-network. Given the image feature maps F_t and the window B it first performs region adaptive max pooling of the features inside B that yields a fixed size feature ($7 \times 7 \times 512$). Then it forwards this feature through two linear+ReLU hidden layers of 4096 channels each (fc_6 and fc_7 layers of VGG16) and a final linear+sigmoid layer with one output that corresponds to the probability p_{obj} of the box B tightly enclosing an object. During training the hidden layers are followed by Dropout units with dropout probability $p = 0.5$.

2.3 Training procedure

Training loss: During training the following multi-task loss is optimized:

$$\underbrace{\frac{1}{N^L} \sum_{k=1}^{N^L} L_{loc}(\theta | B_k, T_k, I_k)}_{\text{localization task loss}} + \underbrace{\frac{1}{N^O} \sum_{k=1}^{N^O} L_{obj}(\theta | B_k, y_k, I_k)}_{\text{objectness scoring task loss}}, \quad (1)$$

where θ are the learnable network parameters, $\{B_k, T_k, I_k\}_{k=1}^{N^L}$ are N^L training triplets for learning the localization task and $\{B_k, y_k, I_k\}_{k=1}^{N^O}$ are N^O training triplets for learning the objectness scoring task. Each training triple $\{B, T, I\}$ of the localization task includes the image I , the box B and the target localization probability vectors $T = \{T_x, T_y\}$. If (B_t^*, B_t^*)

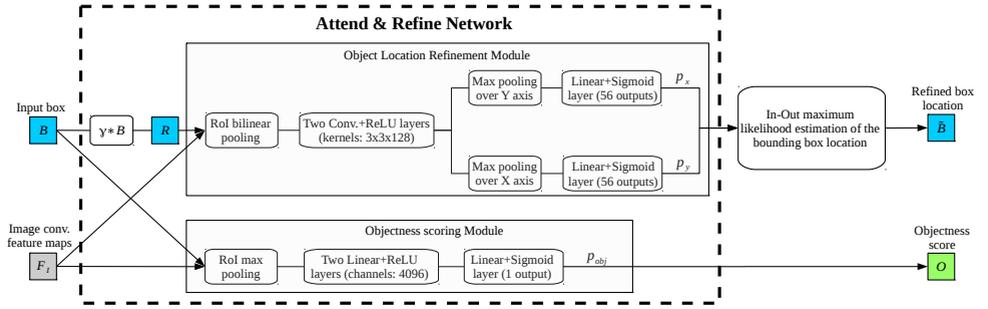


Figure 5: **Attend & Refine Network architecture.** The *Attend & Refine Network* is the box-wise part of the *AttractionNet* architecture. In this figure we depict the work-flow for a single input box B . Specifically, given an input box B and the image convolutional feature maps F_I , the *Attend & Refine Network* yields (1) the in-out location probability vectors, p_x and p_y , (using its object location refinement sub-network) and (2) the objectness scalar probability p_{obj} (using its objectness scoring sub-network). Given the *in-out* probabilities, p_x and p_y , the object location inference is formulated as a simple maximum likelihood estimation problem that results on the refined bounding box coordinates \hat{B} .

and (B_r^*, B_b^*) are the top-left and bottom-right coordinates of the target box B^* then the target probability vectors $T_x = \{T_{x,i}\}_{i=1}^M$ and $T_y = \{T_{y,i}\}_{i=1}^M$ are defined as:

$$T_{x,i} = \begin{cases} 1, & \text{if } B_l^* \leq i \leq B_r^* \\ 0, & \text{otherwise} \end{cases} \text{ and } T_{y,i} = \begin{cases} 1, & \text{if } B_t^* \leq i \leq B_b^* \\ 0, & \text{otherwise} \end{cases}, \forall i \in \{1, \dots, M\} \quad (2)$$

The loss $L_{loc}(\theta|B, T, I)$ of this triplet is the sum of binary logistic regression losses:

$$\frac{1}{2M} \sum_{a \in \{x,y\}} \sum_{i=1}^M T_{a,i} \log(p_{a,i}) + (1 - T_{a,i}) \log(1 - p_{a,i}), \quad (3)$$

where p_a are the output probability vectors of the localization module for the image I , the box B and the network parameters θ . The training triplet $\{B, y, I\}$ for the objectness scoring task includes the image I , the box B and the target value $y \in \{0, 1\}$ of whether the box B contains an object (positive triplet with $y = 1$) or not (negative triplet with $y = 0$). The loss $L_{obj}(\theta|B, y, I)$ of this triplet is the binary logistic regression loss $y \log(p_{obj}) + (1 - y) \log(1 - p_{obj})$, where p_{obj} is the objectness probability for the image I , the box B and the network parameters θ .

Creating training triplets: In order to create the localization and objectness training triplets of one image we first artificially create a pool of boxes that our algorithm is likely to see during test time. Hence we start by generating seed boxes (as the test time algorithm) and for each of them we predict the bounding boxes of the ground truth objects that are closest to them using an ideal object location refinement module. This step is repeated one more time using the previous ideal predictions as input. Because of the finite search area of the search region R the predicted boxes will not necessarily coincide with the ground truth bounding boxes. Furthermore, to account for prediction errors during test time, we repeat the above process by jittering this time the output probability vectors of the ideal location refinement module with 20% noise. Finally, we merge all the generated boxes (starting from the seed ones) to a single pool. Given this pool, the positive training boxes in the objectness localization task are those that their *IoU* with any ground truth object is at least 0.5 and the

negative training boxes are those that their maximum IoU with any ground truth object is less than 0.4. For the localization task we use as training boxes those that their IoU with any ground truth object is at least 0.5.

Optimization: To minimize the objective we use stochastic gradient descent (SGD) optimization with an image-centric strategy for sampling training triplets. Specifically, in each mini-batch we first sample 4 images and then for each image we sample 64 training triplets for the objectness scoring task (50% are positive and 50% are negative) and 32 training triplets for the localization task. The momentum is set to 0.9 and the learning schedule includes training for 320k iterations with a learning rate of $l_r = 0.001$ and then for another 260k iterations with $l_r = 0.0001$. The training time is around 7 days (although we observed that we could have stopped training on the 5th day with insignificant loss in performance).

Scale and aspect ratio jittering: During test time our model is fed with a single image scaled such that its shortest dimension to be 1000 pixels or its longest dimension to not exceed the 1400 pixels. However, during training each image is randomly resized such that its shortest dimension to be one of the following number of pixels $\{300 : 50 : 1000\}$ (using Matlab notation) taking care, however, the longest dimension to not exceed 1000 pixels. Also, with probability 0.5 we jitter the aspect ratio of the image by altering the image dimensions from $W \times H$ to $(\alpha W) \times H$ or $W \times (\alpha H)$ where the value of α is uniformly sampled from $2^{-2} : 2^{1.0}$ (Matlab notation). We observed that this type of data augmentation gives a slight improvement on the results.

3 Experimental results

In this section we perform an exhaustive evaluation of our box proposal generation approach, which we call *AttractionNet*, under various test scenarios. Specifically, we first evaluate our approach with respect to its object localization performance by comparing it with other competing methods and we also provide an ablation study of its main novel components in §3.1. Then, we study its ability to generalize to unseen categories in §3.2, we evaluate it in the context of the object detection task in §3.3 and finally, we provide qualitative results in §3.4.

Training set: In order to train our *AttractionNet* model we use the training set of MS COCO [28] detection benchmark dataset that includes 80k images and it is labelled with 80 different object categories. Note that the MSCOCO dataset is an ideal candidate for training our box proposal model since: (1) it is labelled with a descent number of different object categories and (2) it includes images captured from complex real-life scenes containing common objects in their natural context. The aforementioned training set properties are desirable for achieving good performance on "difficult" test images (a.k.a. images in the wild) and generalizing to "unseen" during training object categories.

Implementation details: In the active box proposal algorithm we use 10k seed boxes generated with a similar to Cracking Bing [43] technique³. To reduce the computational cost of our algorithm, after the first repetition we only keep the top 2k scored boxes and we continue with this number of candidate box proposals for four more extra iterations. In

³We use seed boxes of 3 aspect ratios, 1 : 2, 2 : 1 and 1 : 1, and 9 different sizes of the smallest seed box dimension $\{16, 32, 50, 72, 96, 128, 192, 256, 384\}$.

the non-maximum-suppression [10] (NMS) step the optimal IoU threshold (in terms of the achieved AR) depends on the desired number of box-proposals. For example, for 10, 100, 1000 and 2000 proposals the optimal IoU thresholds are 0.55, 0.75, 0.90 and 0.95 respectively (note that the aforementioned IoU thresholds were cross validated on a set different from the one used for evaluation). For practical purposes and in order to have a unified NMS process, we first apply NMS with the IoU threshold equal to 0.95 and get the top 2000 box proposals, and then follow a multi-threshold NMS strategy that re-orders this set of 2000 boxes such that for any given number K , the top K box proposals in the set better cover (in terms of achieved AR) the objects in the image (see appendix A).

3.1 Object box proposal generation evaluation

Here we evaluate our *AttractionNet* method in the end task of box proposal generation. For that purpose, we test it on the first $5k$ images of the COCO validation set and the PASCAL [10] VOC2007 test set (that also includes around $5k$ images).

Evaluation Metrics: As evaluation metric we use the average recall (AR) which, for a fixed number of box proposals, averages the recall of the localized ground truth objects for several Intersection over Union (IoU) thresholds in the range $.5:.05:.95$ (Matlab notation). The average recall metric has been proposed from Hosang *et al.* [11, 12] where in their work they demonstrated that it correlates well with the average precision performance of box proposal based object detection systems. In our case, in order to evaluate our method we report the AR results for 10, 100 and 1000 box proposals using the notation $AR@10$, $AR@100$ and $AR@1000$ respectively. Also, in the case of 100 box proposals we also report the AR of the small ($\alpha < 32^2$), medium ($32^2 \leq \alpha \leq 96^2$) and large ($\alpha > 96^2$) sized objects using the notation $AR@100\text{-Small}$, $AR@100\text{-Medium}$ and $AR@100\text{-Large}$ respectively, where α is the area of the object. For extracting those measurements we use the COCO API (<https://github.com/pdollar/coco>).

3.1.1 Average recall evaluation

In Table 1 we report the average recall (AR) metrics of our method as well as of other competing methods in the COCO validation set. We observe that the average recall performance achieved by our method exceeds all the previous work in all the AR metrics by a significant margin (around 10 absolute points in the percentage scale). Similar gains are also observed in Table 2 where we report the average recall results of our methods in the PASCAL VOC2007 test set. Furthermore, in Figure 6 we provide for our method the recall as a function of the IoU overlap of the localized ground truth objects. We see that the recall decreases relatively slowly as we increase the IoU from 0.5 to 0.75 while for IoU above 0.85 the decrease is faster.

Comparison with previous state-of-the-art. In Figure 7 we compare the box proposals generated from our *AttractionNet* model (*Ours* entry) against those generated from the previous state-of-the-art [13] (entries *SharpMask*, *SharpMaskZoom* and *SharpMaskZoom²*) w.r.t. the recall versus IoU trade-off and average recall versus proposals number trade-off that they achieve. Also, in Table 1 we report the AR results both for our method and for the

Method	AR@10	AR@100	AR@1000	AR@100-Small	AR@100-Medium	AR@100-Large
EdgeBoxes [14]	0.074	0.178	0.338	0.015	0.134	0.502
Geodesic [14]	0.040	0.180	0.359	-	-	-
Selective Search [14]	0.052	0.163	0.357	0.012	0.0132	0.466
MCG [14]	0.101	0.246	0.398	0.008	0.119	0.530
DeepMask [14]	0.153	0.313	0.446	-	-	-
DeepMaskZoom [14]	0.150	0.326	0.482	-	-	-
Co-Obj [14]	0.189	0.366	0.492	0.107	0.449	0.686
SharpMask [14]	0.192	0.362	0.483	0.060	0.510	0.665
SharpMaskZoom [14]	0.192	0.390	0.532	0.149	0.507	0.630
SharpMaskZoom ² [14]	0.178	0.391	0.555	0.221	0.454	0.588
AttractioNet (Ours)	0.328	0.533	0.662	0.315	0.622	0.777

Table 1: Average Recall results on the first 5k images of COCO validation set.

Method	AR@10	AR@100	AR@1000	AR@100-Small	AR@100-Medium	AR@100-Large
EdgeBoxes [14]	0.203	0.407	0.601	0.035	0.159	0.559
Geodesic [14]	0.121	0.364	0.596	-	-	-
Selective Search [14]	0.085	0.347	0.618	0.017	0.134	0.364
MCG [14]	0.232	0.462	0.634	0.073	0.228	0.618
DeepMask [14]	0.337	0.561	0.690	-	-	-
Best of Co-Obj [14]	0.430	0.602	0.745	0.453	0.517	0.654
AttractioNet (Ours)	0.554	0.744	0.859	0.562	0.670	0.794

Table 2: Average Recall results on the PASCAL VOC2007 test set.

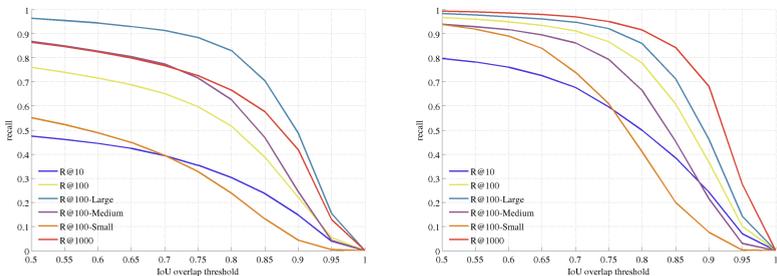


Figure 6: Recall versus IoU overlap plots of our *AttractioNet* approach under different test cases: 10 proposals ($R@10$), 100 proposals ($R@100$), 1000 proposals ($R@1000$), 100 proposals and small sized objects ($R@100$ -Small), 100 proposals and medium sized objects ($R@100$ -Medium) and 100 proposals and large sized objects ($R@100$ -Large). (Left) Results in the first 5k images of COCO validation set. (Right) Results in the PASCAL VOC2007 test set.

Box refinement	Active box generation	# attended boxes	AR@10	AR@100	AR@1000	AR@100-Small	AR@100-Medium	AR@100-Large
		18k	0.147	0.260	0.326	0.122	0.317	0.412
✓		18k	0.298	0.491	0.622	0.281	0.583	0.717
✓	✓	18k	0.328	0.533	0.662	0.315	0.622	0.777

Table 3: Ablation study of our *AttractioNet* box proposal system. In the first row we simply apply the objectness scoring module on a set of 18k seed boxes. In the second row we apply on the same set of 18k seed boxes both the objectness scoring module and the box refinement module. In the last row we utilize our full active box generation strategy that in total attends 18k boxes of which 10k are seed boxes and the rest 8k boxes are actively generated. The reported results are from the first 5k images of COCO validation set.

SharpMask entries. We observe that the model proposed in our work has clearly superior performance over the SharpMask entries under all test cases.

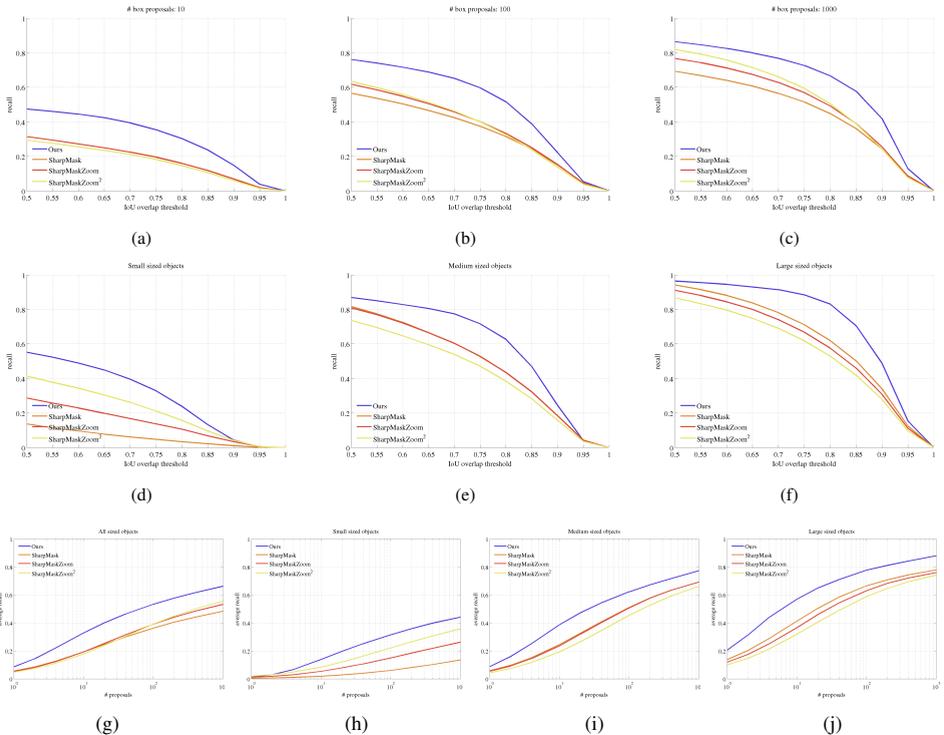


Figure 7: Comparison with previous state-of-the-art. Comparison of our *AttractionNet* box proposal model (*Ours* entry) against the previous state-of-the-art [6] (*SharpMask*, *SharpMaskZoom* and *SharpMaskZoom*² entries) w.r.t. the recall versus IoU trade-off and average recall versus proposals number trade-off that they achieve under various test scenarios. Specifically, the sub-figures (a), (b) and (c) plot the recall as a function of the IoU threshold for 10, 100 and 1000 box proposals respectively and the sub-figures (d), (e) and (f) plot the recall as a function of the IoU threshold for 100 box proposals and with respect to the small, medium and large sized objects correspondingly. Also, the sub-figures (g), (h), (i) and (j) plot the average recall as a function of the proposals number for all the objects regardless of their size as well as for the small, medium and large sized objects respectively. The reported results are from the first 5k images of the COCO validation set.

3.1.2 Ablation study

Here we perform an ablation study of our two key ideas for improving the state-of-the-art on the bounding box proposal generation task:

Object location refinement module. In order to assess the importance of our object location refinement module we evaluated two test cases for generating box proposals: (1) simply applying the objectness scoring module on a set of 18k seed boxes (first row of Table 3) and (2) applying both the objectness scoring module and the object location refinement module on the same set of 18k seed boxes (second row of Table 3). Note that in none of them is the active box generation strategy being used. The average recall results of those two test cases are reported in the first two rows of Table 3. We observe that without the object location refinement module the average recall performance of the box proposal system is very poor. In contrast, the average recall performance of the test case that involves the object

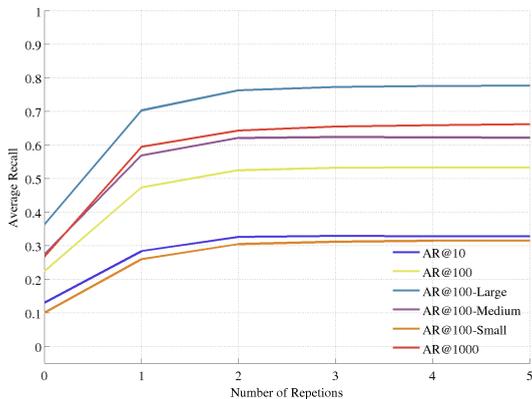


Figure 8: Average recall versus the repetitions number of the active box proposal generation algorithm in the COCO validation set. Note that 0 repetitions is the scenario of simply applying the objectness module on the seed boxes.

location refinement module but not the active box generation strategy is already better than the previous state-of-the-art as reported in Table 1, which demonstrates the very good localization accuracy of our category agnostic location refinement module.

Active box generation strategy. Our active box generation strategy, which we call *Attend Refine Repeat* algorithm, attends in total $18k$ boxes before it outputs the final list of box proposals. Specifically, it attends $10k$ seed boxes in the first repetition of the algorithm and $2k$ actively generated boxes in each of the following four repetitions. A crucial question is whether actively generating those extra $8k$ boxes is really essential in the task or we could achieve the same average recall performance by directly attending $18k$ seed boxes and without continuing on the active box generation stage. We evaluated such test case and we report the average recall results in Table 3 (see rows 2 and 3). We observe that employing the active box generation strategy (3rd row in Table 3) offers a significant boost in the average recall performance (between 3 and 6 absolute points in the percentage scale) thus proving its importance on yielding well localized bounding box proposals. Also, in the right side of Figure 8 we plot the average recall metrics as a function of the repetitions number of our active box generation strategy. We observe that the average recall measurements are increased as we increase the repetitions number and that the increase is more steep on the first repetitions of the algorithm while it starts to converge after the 4th repetition.

3.1.3 Run time

In the current work we did not focus on providing an optimized implementation of our approach. There is room for significantly improving computational efficiency. For instance, just by using SVD decomposition on the fully connected layers of the objectness module at post-training time (similar to Fast-RCNN [14]) and early stopping a sequence of bounding box location refinements in the case it has already converged⁴, the runtime drops from 4.0 seconds to 1.63 seconds without losing almost no accuracy (see Table 4). There are also several other possibilities that we have not yet explored such as tuning the number of feature channels and/or network layers of the CNN architecture (similar to the DeepBox [27] and

⁴ A sequence of bounding box refinements is considered that it has converged when the IoU between the two lastly predicted boxes in the sequence is greater than 0.9.

Method	Run time	AR@10	AR@100	AR@1000	AR@100-Small	AR@100-Medium	AR@100-Large
COCO validation set							
AttractionNet (Ours)	4.00 sec	0.328	0.533	0.662	0.315	0.622	0.777
AttractionNet (Ours, fast version)	1.63 sec	0.326	0.532	0.660	0.317	0.621	0.771
VOC2007 test set							
AttractionNet (Ours)	4.00 sec	0.554	0.744	0.859	0.562	0.670	0.794
AttractionNet (Ours, fast version)	1.63 sec	0.547	0.740	0.848	0.575	0.666	0.788

Table 4: **Run time of our approach on a GTX Titan X GPU.** The reported results are from the first 5k images of COCO validation set and the PASCAL VOC2007 test set.

Method	All categories			Seen categories			Unseen categories		
	AR@10	AR@100	AR@1000	AR@10	AR@100	AR@1000	AR@10	AR@100	AR@1000
AttractionNet (Ours)	0.412	0.618	0.748	0.474	0.671	0.789	0.299	0.521	0.673
EdgeBoxes [14]	0.182	0.377	0.550	0.194	0.396	0.566	0.160	0.344	0.519
Selective Search [15]	0.132	0.358	0.562	0.143	0.372	0.568	0.111	0.332	0.551
MCG [16]	0.219	0.428	0.603	0.228	0.447	0.623	0.205	0.395	0.568

Table 5: **Generalization to unseen categories: from COCO to ImageNet.** In this table we report average recall results on the ImageNet [17] ILSVRC2013 detection task validation set that includes around 20k images and it is labelled with 200 object categories. *Seen categories* are the set of object categories that our COCO trained *AttractionNet* model "saw" during training. In contrast, *unseen categories* is the set of object categories that were not present in the training set of our *AttractionNet* model.

the SharpMask [17] approaches).

In the remainder of this section we will use the fast version of our *AttractionNet* approach in order to provide experimental results.

3.2 Generalization to unseen categories

So far we have evaluated our *AttractionNet* approach — in the end task of object box proposal generation — on the COCO validation set and the PASCAL VOC2007 test set that are labelled with the same or a subset of the object categories "seen" in the training set. In order to assess the *AttractionNet*'s capability to generalize to "unseen" categories, as it is suggested by Chavali *et al.* [9], we evaluate our *AttractionNet* model on two extra datasets that are labelled with object categories that are not present in its training set ("unseen" object categories).

From COCO to ImageNet [17]. Here we evaluate our COCO trained *AttractionNet* box proposal model on the ImageNet [17] ILSVRC2013 detection task validation set that is labelled with 200 different object categories and we report average recall results in Table 5. Note that among the 200 categories of ImageNet detection task, 60 of them, as we identified, are also present in the *AttractionNet*'s training set (see Appendix C). Thus, for a better insight on the generalization capabilities of *AttractionNet*, we divided the ImageNet detection task categories on two groups, the "seen" by *AttractionNet* categories and the "unseen" categories, and we report the average recall results separately for those two groups of object categories in Table 5. For comparison purposes we also report the average recall performance of a few indicative other box proposal methods that their code is publicly available. We observe that, despite the performance difference of our approach between the "seen" and the "unseen" object categories (which is to be expected), its average recall performance on the "unseen" categories is still quite high and significantly better than the other box proposal methods. Note that even the non-learning based approaches of Selective Search and EdgeBoxes exhibit a performance drop on the "unseen" by *AttractionNet* group of object categories, which we assume is because this group contains more intrinsically difficult to discover objects.

Method	AR@10	AR@100	AR@1000	AR@100-Small	AR@100-Medium	AR@100-Large
AttractionNet (Ours)	0.159	0.389	0.579	0.205	0.419	0.498
EdgeBoxes [14]	0.049	0.160	0.362	0.020	0.131	0.332
Selective Search [15]	0.024	0.143	0.422	0.008	0.085	0.362
MCG [16]	0.078	0.237	0.441	0.045	0.195	0.476

Table 6: **Generalization to unseen categories: from COCO to NYU-Depth V2 dataset.** In this table we report average recall results on the 1449 labelled images of the NYU-Depth V2 dataset [17]. Note that the NYU-Depth V2 dataset is densely labelled with more than 800 different categories.

From COCO to NYU Depth dataset [17]. The NYU Depth V2 dataset [17] provides 1449 images (recorded from indoor scenes) that are densely pixel-wise annotated with 864 different categories. We used the available instance-wise segmentations to create ground truth bounding boxes and we tested our COCO trained *AttractionNet* model on them (see Table 6). Note that among the 864 available pixel categories, a few of them are "stuff" categories (e.g. wall, floor, ceiling or stairs) or in general non-object pixel categories that our object box proposal method should by definition not recall. Thus, during the process of creating the ground truth bounding boxes, those non-object pixel segmentation annotations were excluded (see Appendix D). In Table 6 we report the average recall results of our *AttractionNet* method as well as of a few other indicative methods that their code is publicly available. We again observe that our method surpasses all other approaches by a significant margin. Furthermore, in this case the superiority of our approach is more evident on the average recall of the small and medium sized objects.

To conclude we argue that our learning based *AttractionNet* approach exhibits good generalization behaviour. Specifically, its average recall performance on the "unseen" object categories remains very high and is also much better than other competing approaches, including both learning-based approaches such as the MSG and hand-engineered ones such as the Selective Search or the EdgeBoxes methods. A performance drop is still observed while going from "seen" to "unseen" categories, but this is something to be expected given that any machine learning algorithm will always exhibit a certain performance drop while going from "seen" to "unseen" data (i.e. training set accuracy versus test set accuracy).

3.3 AttractionNet box proposals evaluation in the context of the object detection task

Here we evaluate our *AttractionNet* box proposals in the context of the object detection task by training and testing a box proposal based object detection system on them (specifically we use the fast version of *AttractionNet* that is described in section 3.1.3).

Detection system. Our box proposal based object detection system consists of a Fast-RCNN [18] category-specific recognition module and a LocNet Combined ML [19] category-specific bounding box refinement module (see Appendix B for more details). As post-processing we use a non-max-suppression step (with IoU threshold of 0.35) that is enhanced with the box voting technique described in the MR-CNN system [13] (with IoU threshold of 0.75). Note that we did not include iterative object localization as in the LocNet [19] or MR-CNN [13] papers, since our bounding box proposals are already very well localized and we did not get any significant improvement from running the detection system for extra iterations. Using the same trained model we provide results for two test cases: (1) using a single scale of 600 pixels during test time and (2) using two scales of 500 and 1000 pixels

during test time.

Detection evaluation setting. The detection evaluation metrics that we use are the average precision (AP) for the IoU thresholds of 0.50 (AP@0.50), 0.75 (AP@0.75) and the COCO style of average precision (AP@0.50 : 0.95) that averages the traditional AP over several IoU thresholds between 0.50 and 0.95. Also, we report the COCO style of average precision with respect to the small (AP@Small), medium (AP@Medium) and large (AP@Large) sized objects. We perform the evaluation on 5k images of COCO 2014 validation set and we provide final results on the COCO 2015 test-dev set.

Detection results. In Figure 9 we provide plots of the achieved average precision (AP) as a function of the used box proposals number and in Table 7 we provide the average precision results for 10, 100, 1000 and 2000 box proposals. We observe that in all cases, the average precision performance of the detection system seems to converge after the 200 box proposals. Furthermore, for single scale test case our best COCO-style average precision is 0.320 and for the two scales test case our best COCO-style average precision is 0.337. By including horizontal image flipping augmentation during test time our COCO-style average precision performance is increased to 0.343. Finally, in Table 8 we provide the average precision performance in the COCO test-dev 2015 set where we achieve a COCO-style AP of 0.341. By comparing with the average precision performance of the other competing methods, we observe that:

- Comparing with the other VGG16-Net based object detection systems (ION [9] and MultiPath [22] systems), our detection system achieves the highest COCO-style average precision with its main novelties w.r.t. the Fast R-CNN [15] baseline being (1) the use of the *AttractionNet* box proposals that are introduced in this paper and (2) the LocNet [24] category specific object location refinement technique that replaces the bounding box regression step.
- Comparing with the ION [9] detection system, which is also VGG16-Net based, our approach is better on the COCO-style AP metric (that favours good object localization) while theirs is better on the typical AP@0.50 metric. We hypothesize that this is due to the fact that our approach targets to mainly improve the localization aspect of object detection by improving the box proposal generation step while theirs the recognition aspect of object detection. The above observation suggests that many of the novelties introduced on the ION [9] and MultiPath [22] systems w.r.t. object detection could be orthogonal to our box proposal generation work.
- The achieved average precision performance of our VGG16-Net based detection system is close to the state-of-the-art *ResNet-101 based* Faster R-CNN+++ detection system [20] that exploits the recent successes in deep representation learning introduced — under the name Deep Residual Networks — in the same work by He *et al.* [20]. Presumably, our overall detection system could also be benefited by being based on the Deep Residual Networks [20] or the more recent wider variant called Wide Residual Networks [44].
- Finally, our detection system has the highest average precision performance w.r.t. the small sized objects, which is a challenging problem, surpassing by a healthy margin even the ResNet-101 based Faster R-CNN+++ detection system [20]. This is thanks to the high average recall performance of our box proposal method on the small sized objects.

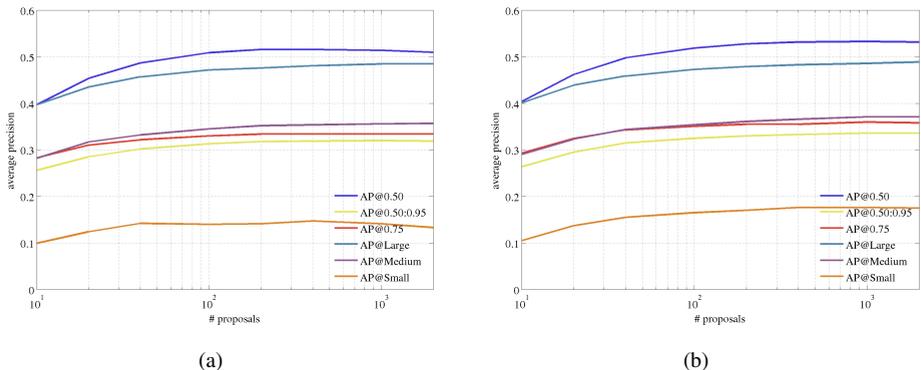


Figure 9: **Detection results: Average precision versus *AttractionNet* box proposals number.** (a) During test time a single scale of 600 pixels is being used. (b) During test time two scales of 500 and 1000 pixels are being used. The reported results are from 5k images of COCO validation set.

Test scale(s)	# proposals	AP@0.50	AP@0.75	AP@0.50:0.95	AP@Small	AP@Medium	AP@Large
600px	10	0.397	0.283	0.256	0.099	0.282	0.397
600px	100	0.509	0.330	0.313	0.140	0.345	0.472
600px	1000	0.514	0.334	0.320	0.141	0.356	0.485
600px	2000	0.510	0.334	0.319	0.133	0.357	0.485
500px, 1000px	10	0.404	0.293	0.264	0.105	0.290	0.401
500px, 1000px	100	0.519	0.351	0.325	0.165	0.354	0.473
500px, 1000px	1000	0.533	0.360	0.336	0.176	0.371	0.486
500px, 1000px	2000	0.532	0.358	0.336	0.175	0.371	0.489
500px, 1000px ★	2000	0.540	0.364	0.343	0.184	0.382	0.491

Table 7: **Detection results: Average precision performance using *AttractionNet* box proposals.** The reported results are from 5k images of COCO validation set. The last entry with the ★ symbol uses horizontal image flipping augmentation during test time.

3.4 Qualitative results

In Figure 10 we provide qualitative results of our *AttractionNet* box proposal approach on images coming from the COCO validation set. Note that our approach manages to recall most of the objects in an image, even in the case that the depicted scene is crowded with multiple objects that heavily overlap with each other.

4 Conclusions

In our work we propose a bounding box proposals generation method, which we call *AttractionNet*, whose key elements are a strategy for actively searching of bounding boxes in the promising image areas and a powerful object location refinement module that extends the recently introduced LocNet [14] model on localizing objects agnostic to their category. We extensively evaluate our method on several image datasets (i.e. COCO, PASCAL, ImageNet detection and NYU-Depth V2 datasets) demonstrating in all cases average recall results that surpass the previous state-of-the-art by a significant margin while also providing strong empirical evidence about the generalization ability of our approach w.r.t. unseen categories. Even more, we show the significance of our *AttractionNet* approach in the object detection task by coupling it with a VGG16-Net based detector and thus managing to surpass the detection performance of all other VGG16-Net based detectors while even being on par with a heavily tuned ResNet-101 based detector. We note that, apart from object detection, there



Figure 10: **Qualitative results in COCO.** The blue rectangles are the box proposals generated by our approach that best localize (in terms of IoU) the ground truth boxes. The red rectangles are the ground truth bounding boxes that were not discovered by our box proposal approach (their IoU with any box proposal is less than 0.5). Note that not all the object instances on the images are annotated.

Method	Base CNN	AP@0.50	AP@0.75	AP@0.50:0.95	AP@Small	AP@Medium	AP@Large
<i>AttractionNet</i> based detection system (Ours)	VGG16-Net [59]	0.537	0.363	0.341	0.175	0.365	0.469
ION [6]	VGG16-Net [59]	0.557	0.346	0.331	0.145	0.352	0.472
MultiPath [42]	VGG16-Net [59]	-	-	0.315	-	-	-
Faster R-CNN+++ [24]	ResNet-101 [59]	0.557	-	0.349	0.156	0.387	0.509

Table 8: Detection results in COCO test-dev 2015 set. In this table we report the average precision performance of our *AttractionNet* box proposals based detection system that uses 2000 proposals and two test scales of 500 and 1000 pixels. Note that: (1) all methods in this table (including ours) use horizontal image flipping augmentation during test time, (2) the ION [6] and MultiPath [42] detection systems use a single test scale of 600 and 800 pixels respectively while the Faster R-CNN+++ entry uses the scales {200, 400, 600, 800, 1000}, (3) apart from the ResNet-101 based Faster R-CNN+++ [24] entry, all the other methods are based on the VGG16-Network [59], (4) the reported results of all the competing methods are from the single model versions of their systems (and not the model ensemble versions) and (5) the reported results of the MultiPath system are coming from 5k images of the COCO validation set (however, we expect the AR metrics on the test-dev set to be roughly similar).

exist several other vision tasks, such as exemplar 2D-3D detection [60], visual semantic role labelling [17], caption generation [24] or visual question answering [62], for which a box proposal generation step can be employed. We are thus confident that our *AttractionNet* approach could have a significant value with respect to many other important applications as well.

Acknowledgements

This work was supported by the ANR SEMAPOLIS project. We would like to thank Pedro O. Pinheiro for help with the experimental results and Sergey Zagoruyko for helpful discussions. We would also like to thank the authors of SharpMask [62] (Pedro O. Pinheiro, Tsung-Yi Lin, Ronan Collobert and Piotr Dollar) for providing us with its box proposals.

References

- [1] Bogdan Alexe, Thomas Deselaers, and Vittorio Ferrari. Measuring the objectness of image windows. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2012.
- [2] P. Arbeláez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In *Computer Vision and Pattern Recognition*, 2014.
- [3] Sean Bell, C Lawrence Zitnick, Kavita Bala, and Ross Girshick. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. *arXiv preprint arXiv:1512.04143*, 2015.
- [4] Neelima Chavali, Harsh Agrawal, Aroma Mahendru, and Dhruv Batra. Object-proposal evaluation protocol is ‘gameable’. *arXiv preprint arXiv:1505.05836*, 2015.
- [5] Xiaozhi Chen, Huimin Ma, Xiang Wang, and Zhichen Zhao. Improving object proposals with multi-thresholding straddling expansion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2587–2595, 2015.

- [6] Ming-Ming Cheng, Ziming Zhang, Wen-Yan Lin, and Philip Torr. Bing: Binarized normed gradients for objectness estimation at 300fps. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, 2014.
- [7] Ramazan Gokberk Cinbis, Jakob Verbeek, and Cordelia Schmid. Weakly supervised object localization with multi-fold multiple instance learning. *arXiv preprint arXiv:1503.00949*, 2015.
- [8] Jifeng Dai, Kaiming He, Yi Li, Shaoqing Ren, and Jian Sun. Instance-sensitive fully convolutional networks. *CoRR*, abs/1603.08678.
- [9] Jifeng Dai, Kaiming He, and Jian Sun. Instance-aware semantic segmentation via multi-task network cascades. *arXiv preprint arXiv:1512.04412*, 2015.
- [10] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 2010.
- [11] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2010.
- [12] Amir Ghodrati, Ali Diba, Marco Pedersoli, Tinne Tuytelaars, and Luc Van Gool. Deep-proposal: Hunting objects by cascading deep convolutional layers. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2578–2586, 2015.
- [13] Spyros Gidaris and Nikos Komodakis. Object detection via a multi-region & semantic segmentation-aware cnn model. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [14] Spyros Gidaris and Nikos Komodakis. Locnet: Improving localization accuracy for object detection. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on Computer Vision*, 2016.
- [15] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [16] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, 2014.
- [17] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015.
- [18] Zeeshan Hayder, Xuming He, and Mathieu Salzmann. Learning to co-generate object proposals with a deep structured network. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, number EPFL-CONF-217984, 2016.
- [19] K He, X Zhang, S Ren, and J Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2015.

- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [21] J. Hosang, R. Benenson, and B. Schiele. How good are detection proposals, really? In *BMVC*, 2014.
- [22] Jan Hosang, Rodrigo Benenson, Piotr Dollár, and Bernt Schiele. What makes for effective detection proposals? *PAMI*, 2015.
- [23] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Denscap: Fully convolutional localization networks for dense captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [24] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- [25] Philipp Krähenbühl and Vladlen Koltun. Geodesic object proposals. In *Computer Vision–ECCV 2014*. Springer, 2014.
- [26] Philipp Krähenbühl and Vladlen Koltun. Learning to propose objects. In *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*. IEEE, 2015.
- [27] Weicheng Kuo, Bharath Hariharan, and Jitendra Malik. Deepbox: Learning objectness with convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2479–2487, 2015.
- [28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014*, 2014.
- [29] Yongxi Lu, Tara Javidi, and Svetlana Lazebnik. Adaptive object detection using adjacency and zoom prediction. *CoRR*, abs/1512.07711.
- [30] Francisco Massa, Bryan Russell, and Mathieu Aubry. Deep exemplar 2d-3d detection by adapting from real to rendered views. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on Computer Vision*, 2016.
- [31] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, 2012.
- [32] Pedro H. O. Pinheiro, Tsung-Yi Lin, Ronan Collobert, and Piotr Dollár. Learning to refine object segments. 2016. URL <http://arxiv.org/abs/1603.08695>.
- [33] Pedro O Pinheiro, Ronan Collobert, and Piotr Dollar. Learning to segment object candidates. In *Advances in Neural Information Processing Systems*, 2015.
- [34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, 2015.

- [35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 2015.
- [36] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- [37] Kevin J Shih, Saurabh Singh, and Derek Hoiem. Where to look: Focus regions for visual question answering. *arXiv preprint arXiv:1511.07394*, 2015.
- [38] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. *arXiv preprint arXiv:1604.03540*, 2016.
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [40] Koen EA Van de Sande, Jasper RR Uijlings, Theo Gevers, and Arnold WM Smeulders. Segmentation as selective search for object recognition. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, 2011.
- [41] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *CoRR*, abs/1605.07146, 2016.
- [42] Sergey Zagoruyko, Adam Lerer, Tsung-Yi Lin, Pedro O Pinheiro, Sam Gross, Soumith Chintala, and Piotr Dollár. A multipath network for object detection. *arXiv preprint arXiv:1604.02135*, 2016.
- [43] Qiyang Zhao, Zhibin Liu, and Baolin Yin. Cracking bing and beyond. In *Proceedings of the British Machine Vision Conference. BMVA Press*, 2014.
- [44] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *Computer Vision—ECCV 2014*, 2014.

A Multi-threshold non-max-suppression re-ordering

As already described in section 2.1, at the end of our active box proposal generation strategy we include a non-maximum-suppression (\square) (NMS) step that is applied on the set \mathbf{C} of scored candidate box proposals in order then to take the final top K output box proposals (see algorithm 1). However, the optimal IoU threshold (in terms of the achieved AR) for the NMS step depends on the desired number K of output box-proposals. For example, for 10, 100, 1000 and 2000 proposals the optimal IoU thresholds are 0.55, 0.75, 0.90 and 0.95 respectively. Since our plan is to make our box proposal system publicly available, we would like to make its use easier for the end user. For that purpose, we first apply on the set \mathbf{C} of scored candidate box proposals a simple NMS step with IoU threshold equal to 0.95 in order to then get the top 2000 box proposals and then we follow a multi-threshold non-max-suppression technique that re-orders this set of 2000 box proposals such that for any given number K the top K box proposals in the set better cover (in terms of achieved AR) the objects in the image.

Specifically, assume that $\{t_i\}_{i=1}^{N_k}$ are the optimal IoU thresholds for N_k different desired numbers of output box proposals $\{K_i\}_{i=1}^{N_k}$, where both the thresholds and the desired number of box proposals are in ascending order⁵. Our multi-threshold NMS strategy starts by applying on the aforementioned set \mathbf{L} of 2000 box proposals simple single-threshold NMS steps with IoU thresholds $\{t_i\}_{i=1}^{N_k}$ that results on N_k different lists of box proposals $\{\mathbf{L}(t_i)\}_{i=1}^{N_k}$ (note that all the NMS steps are applied on the same list \mathbf{L} and not in consecutive order). Then, starting from the lowest threshold t_1 (which also is the more restrictive one) we take from the list $\mathbf{L}(t_1)$ the top K_1 box proposals and we add them to the set of output box proposals \mathbf{P} . For the next threshold t_2 we get the top $K_2 - |\mathbf{P}|$ box proposals from the set $\{\mathbf{L}(t_2) \setminus \mathbf{P}\}$ and again add them on the set \mathbf{P} . This process continues till the last threshold t_{N_k} at which point the size of the output box proposals set \mathbf{P} is $K_{N_k} = 2000$. Each time $i = 1, \dots, N_k$ we add box proposals on the set \mathbf{P} , their objectness scores are altered according to the formula $\bar{o} = o + (N_k - i)$ (where o and \bar{o} are the initial and after re-ordering objectness scores correspondingly) such that their new objectness scores to correspond to the order at which they are placed in the set \mathbf{P} . Note that this technique does not guarantee an optimal re-ordering of the boxes (in terms of AR), however it works sufficiently well in practice.

B Detection system

In this section we provide further implementation details about the object detection system used in §3.3.

Architecture. Our box proposal based object detection network consists of a Fast-RCNN [15] category-specific recognition module and a LocNet Combined ML [14] category-specific bounding box refinement module that share the same image-wise convolutional layers (conv1_1 till conv5_3 layers of VGG16-Net).

Training. The detection network is trained on the union of the COCO train set that includes around 80k images and on a subset of the COCO validation set that includes around 35k images (the remaining 5k images of COCO validation set are being used for evaluation). For training we use our *AttractionNet* box proposals and we define as positives those that have IoU overlap with any ground truth bounding box at least 0.5 and as negatives the remaining proposals. For training we use SGD where each mini-batch consists of 4 images with 64 box proposals each (256 boxes per mini-batch in total) and the ratio of negative-positive boxes is 3:1. We train the detection network for 500k SGD iterations starting with a learning rate of 0.001 and dropping it to 0.0001 after 320k iterations. We use the same scale and aspect ratio jittering technique that is used on *AttractionNet* and is described in section 2.3.

C Common categories between ImageNet and COCO

In this section we list the ImageNet detection task object categories that we identified to be present also in the COCO dataset. Those are:
airplane, apple, backpack, baseball, banana, bear, bench, bicycle, bird, bowl, bus,

⁵We used the IoU thresholds of $\{0.55, 0.60, 0.65, 0.75, 0.80, 0.85, 0.90, 0.95\}$ for the desired numbers of output box proposals $\{10, 20, 40, 100, 200, 400, 1000, 2000\}$. Those IoU thresholds were cross-validated on a validation set different than this used for the evaluation of our approach.

car, chair, cattle, computer keyboard, computer mouse, cup or mug, dog, domestic cat, digital clock, elephant, horse, hotdog, laptop, microwave, motorcycle, orange, person, pizza, refrigerator, sheep, ski, tie, toaster, traffic light, train, zebra, racket, remote control, sofa, tv or monitor, table, watercraft, washer, water bottle, wine bottle, ladle, flower pot, purse, stove, koala bear, volleyball, hair dryer, soccer ball, rugby ball, croquet ball, basketball, golf ball, ping-pong ball, tennis ball.

D Ignored NUY Depth dataset categories

In this section we list the 12 most frequent non-object categories that we identified on the NUY Depth V2 dataset:

curtain, cabinet, wall, floor, ceiling, room divider, window shelf, stair, counter, window, pipe and column.