

What makes for effective detection proposals?

Jan Hosang¹, Rodrigo Benenson¹, Piotr Dollár², and Bernt Schiele¹

¹Max Planck Institute for Informatics

²Facebook AI Research (FAIR)

Abstract—Current top performing object detectors employ detection proposals to guide the search for objects, thereby avoiding exhaustive sliding window search across images. Despite the popularity and widespread use of detection proposals, it is unclear which trade-offs are made when using them during object detection. We provide an in-depth analysis of twelve proposal methods along with four baselines regarding proposal repeatability, ground truth annotation recall on PASCAL, ImageNet, and MS COCO, and their impact on DPM, R-CNN, and Fast R-CNN detection performance. Our analysis shows that for object detection improving proposal localisation accuracy is as important as improving recall. We introduce a novel metric, the average recall (AR), which rewards both high recall and good localisation and correlates surprisingly well with detection performance. Our findings show common strengths and weaknesses of existing methods, and provide insights and metrics for selecting and tuning proposal methods.

Index Terms—Computer Vision, object detection, detection proposals.



1 INTRODUCTION

UNTIL recently, the most successful approaches to object detection utilised the well known “sliding window” paradigm [1]–[3], in which a computationally efficient classifier tests for object presence in every candidate image window. Sliding window classifiers scale linearly with the number of windows tested, and while single-scale detection requires classifying around $10^4 - 10^5$ windows per image, the number of windows grows by an order of magnitude for multi-scale detection. Modern detection datasets [4]–[6] also require the prediction of object aspect ratio, further increasing the search space to $10^6 - 10^7$ windows per image.

The steady increase in complexity of the core classifiers has led to improved detection quality, but at the cost of significantly increased computation time per window [7]–[11]. One approach for overcoming the tension between computational tractability and high detection quality is through the use of “detection proposals” [12]–[15]. Under the assumption that all objects of interest share common visual properties that distinguish them from the background, one can design or train a method that, given an image, outputs a set of proposal regions that are likely to contain objects. If high object recall can be reached with considerably fewer windows than used by sliding window detectors, significant speed-ups can be achieved, enabling the use of more sophisticated classifiers.

Current top performing object detectors for PASCAL [4] and ImageNet [5] all use detection proposals [7]–[11], [16]. In addition to allowing for use of more sophisticated classifiers, the use of detection proposals alters the data distribution that the classifiers handle. This may also improve detection quality by reducing spurious false positives.

Most papers on generating detection proposals perform fairly limited evaluations, comparing results using only a subset of metrics, datasets, and competing methods. In this work, we aim to revisit existing work on proposals

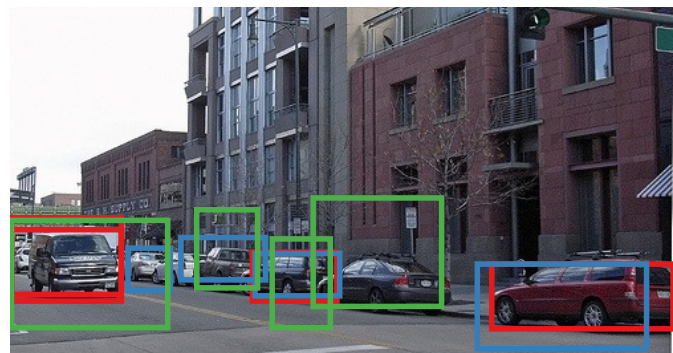


Figure 1: What makes object detection proposals effective?

and compare most publicly available methods in a unified framework. While this requires us to carefully re-examine the metrics and settings for evaluating proposals, it allows us to better understand the benefits and limitations of current methods.

The contributions of this work are as follows:

- In §2 we provide a systematic overview of detection proposal methods and define simple baselines that serve as reference points. We discuss the taxonomy of proposal methods, and describe commonalities and differences of the various approaches.
- In §3 we introduce the notion of proposal repeatability, discuss its relevance when considering proposals for detection, and measure the repeatability of existing methods. The results are somewhat unexpected.
- In §4 we study object recall on the PASCAL VOC 2007 test set [4], and for the first time, over the larger and more diverse ImageNet 2013 [5] and MS COCO 2014 [6] validation sets. The latter allows us to examine possible biases towards PASCAL objects categories. Overall, these experiments are substan-

tially broader in scope than previous work, both in the number of methods evaluated and datasets used.

- In §5 we evaluate the influence of different proposal methods on DPM [3], R-CNN [8], and Fast R-CNN [16] detection performance. Based on our results, we introduce a novel evaluation metric, the average recall (AR). We show that AR is highly correlated with detector performance, more so than previous metrics, and we advocate AR to become the standard metric for evaluating proposals. Our experiments provide the first clear guidelines for selecting and tuning proposal methods for object detection.

All evaluation scripts and method bounding boxes used in this work are publicly available to facilitate the reproduction of our evaluation¹. The results presented in this paper summarise results of over 500 experiments on multiple data sets and required multiple months of CPU time.

An earlier version of this work appeared in [17].

2 DETECTION PROPOSAL METHODS

Detection proposals are similar in spirit to interest point detectors [30], [31]. Interest points allow for focusing attention to the most salient and distinctive locations in an image, greatly reducing computation for subsequent tasks such as classification, retrieval, matching, and detection. Likewise, object proposals considerably reduce computation compared to the dense (sliding window) detection framework by generating candidate proposals that may contain objects. This in turn enables use of expensive classifiers per window [7]–[11].

It is worthwhile noting that interest points were dominant when computing feature descriptors densely was prohibitive. However, with improved algorithmic efficiency and increased computational power, it is now standard practice to use dense feature extraction [32]. The opposite trend has occurred in object detection, where the dense sliding window framework has been overtaken by use of proposals. We aim to understand if detection proposals improve detection accuracy or if their use is strictly necessary for computational reasons. While in this work we focus on the impact of proposals on detection, proposals have applications beyond object detection, as we discuss in §6.

Two general approaches for generating object proposals have emerged: *grouping methods* and *window scoring methods*. These are perhaps best exemplified by the early and well known *SelectiveSearch* [15] and *Objectness* [12] proposal methods. We survey these approaches in §2.1 and §2.2, followed by an overview of alternate approaches in §2.3 and baselines in §2.4. Finally, we consider the connection between proposals and cascades in §2.5 and provide additional method details in §2.6.

The survey that follows is meant to be exhaustive. However, for the purpose of our evaluations, we only consider methods for which source code is available. We cover a diverse set of methods (in terms of quality, speed, and underlying approach). Table 1 gives an overview of the 12 selected methods (plus 4 baselines).² Table 1 also indicates

high level information regarding the output of each method and a qualitative overview of the results of the evaluations performed in the remainder of this paper.

In this paper we concentrate on class-agnostic proposals for single-frame, bounding box detection. For proposal methods that output segmentations instead of bounding boxes, we convert the output to bounding boxes for the purpose of our evaluation. Methods that operate on videos and require temporal information (e.g. [33]) are considered outside the scope of this work.

2.1 Grouping proposal methods

Grouping proposal methods attempt to generate multiple (possibly overlapping) segments that are likely to correspond to objects. The simplest such approach would be to directly use the output of any hierarchical image segmentation algorithm, e.g. Gu et al. [34] use the segmentation produced by gPb [35]. To increase the number of candidate segments, most methods attempt to diversify such hierarchies, e.g. by using multiple low level segmentations [19], [26], [29] or starting with an over-segmentation and randomising the merge process [26]. The decision to merge segments is typically based on a diverse set of cues including superpixel shape, appearance cues, and boundary estimates (typically obtained from [35], [36]).

We classify grouping methods into three types according to how they generate proposals. Broadly speaking, methods generate region proposals by grouping superpixels (SP), often using [37], solving multiple graph cut (GC) problems with diverse seeds, or directly from edge contours (EC), e.g. from [35], [36]. In the method descriptions below the type of each method is marked by SP, GC, or EC accordingly.

We note that while all the grouping approaches have the strength of producing a segmentation mask of the object, we evaluate only the enclosing bounding box proposals.

- **SelectiveSearch**^{†SP} [15], [29] greedily merges superpixels to generate proposals. The method has no learned parameters, instead features and similarity functions for merging superpixels are manually designed. *SelectiveSearch* has been broadly used as the proposal method of choice by many state-of-the-art object detectors, including the R-CNN and Fast R-CNN detectors [8], [16].
- **RandomizedPrim's**^{†SP} [26] uses similar features as *SelectiveSearch*, but introduces a randomised superpixel merging process in which all probabilities have been learned. Speed is substantially improved.
- **Rantalankila**^{†SP} [27] proposes a superpixel merging strategy similar to *SelectiveSearch*, but using different features. In a subsequent stage, the generated segments are used as seeds for solving graph cuts in the spirit of CPMC (see below) to generate more proposals.
- **Chang**^{SP} [38] combines saliency and *Objectness* with a graphical model to merge superpixels into figure/background segmentations.
- **CPMC**^{†GC} [13], [19] avoids initial segmentations and computes graph cuts with several different seeds and unaries directly on pixels. The resulting segments are ranked using a large pool of features.

1. Project page: <http://goo.gl/uMhkAs>

2. We mark the evaluated methods with a '†' in the following listing.

Method	Approach	Outputs Segments	Outputs Score	Control #proposals	Time (sec.)	Repeatability	Recall Results	Detection Results
Bing [18]	Window scoring		✓	✓	0.2	***	*	.
CPMC [19]	Grouping	✓	✓	✓	250	-	**	*
EdgeBoxes [20]	Window scoring		✓	✓	0.3	**	***	***
Endres [21]	Grouping	✓	✓	✓	100	-	***	**
Geodesic [22]	Grouping	✓	✓	✓	1	*	***	**
MCG [23]	Grouping	✓	✓	✓	30	*	***	***
Objectness [24]	Window scoring		✓	✓	3	.	*	.
Rahtu [25]	Window scoring		✓	✓	3	.	.	*
RandomizedPrim's [26]	Grouping	✓		✓	1	*	*	**
Rantalankila [27]	Grouping	✓		✓	10	**	.	**
Rigor [28]	Grouping	✓		✓	10	*	**	**
SelectiveSearch [29]	Grouping	✓	✓	✓	10	**	***	***
Gaussian				✓	0	.	.	*
SlidingWindow				✓	0	***	.	.
Superpixels		✓			1	*	.	.
Uniform				✓	0	.	.	.

Table 1: Comparison of different detection proposal methods. Grey check-marks indicate that the number of proposals is controlled by indirectly adjusting parameters. Repeatability, quality, and detection rankings are provided as rough summary of the experimental results: “-” indicates no data, “.”, “*”, “**”, “***” indicate progressively better results. These guidelines were obtained based on experiments presented in sections §3, §4, and §5, respectively.

- **Endres**^{†GC} [14], [21] builds a hierarchical segmentation from occlusion boundaries and solves graph cuts with different seeds and parameters to generate segments. The proposals are ranked based on a wide range of cues and in a way that encourages diversity.
- **Rigor**^{†GC} [28] is a somewhat improved variant of CPMC that speeds computation considerably by re-using computation across multiple graph-cut problems and using the fast edge detectors from [36], [39].
- **Geodesic**^{†EC} [22] starts from an over-segmentation of the image based on [36]. Classifiers are used to place seeds for a geodesic distance transform. Level sets of each of the distance transforms define the figure/ground segmentations that are the proposals.
- **MCG**^{†EC} [23] introduces a fast algorithm for computing multi-scale hierarchical segmentations building on [36]. Segments are merged based on edge strength and the resulting object hypotheses are ranked using cues such as size, location, shape, and edge strength.

2.2 Window scoring proposal methods

An alternate approach for generating detection proposals is to score each candidate window according to how likely it is to contain an object. Compared to grouping approaches these methods usually only return bounding boxes and tend to be faster. Unless window sampling is performed very densely, this approach typically generates proposals with low localisation accuracy. Some methods counteract this by refining the location of the generated windows.

- **Objectness**[†] [12], [24] is one of the earliest and well known proposal methods. An initial set of proposals is selected from salient locations in an image, these proposals are then scored according to multiple cues including colour, edges, location, size, and the strong “superpixel straddling” cue.
- **Rahtu**[†] [25] begins with a large pool of proposal regions generated from individual superpixels, pairs and triplets of superpixels, and multiple randomly sampled boxes. The scoring strategy used by Ob-

jectness is revisited, and improvements are proposed. [40] adds additional low-level features and highlights the importance of properly tuned non-maximum suppression.

- **Bing**[†] [18] uses a simple linear classifier trained over edge features and applied in a sliding window manner. Using adequate approximations a very fast class agnostic detector is obtained (1 ms/image on CPU). However, it was shown that the classifier has minimal influence and similar performance can be obtained *without* looking at the image [41]. This image independent method is named `CrackingBing`.
- **EdgeBoxes**^{†EC} [20] also starts from a coarse sliding window pattern, but builds on object boundary estimates (obtained via structured decision forests [36], [42]) and adds a subsequent refinement step to improve localisation. No parameters are learned. The authors propose tuning the density of the sliding window pattern and the threshold of the non-maximum suppression to tune the method for different overlap thresholds (see §5).
- **Feng** [43] poses proposal generation as the search for salient image content and introduces new saliency measures, including the ease with which a potential object can be composed from the rest of the image. The sliding window paradigm is used and every location scored according to the saliency cues.
- **Zhang** [44] proposes to train a cascade of ranking SVMs on simple gradient features. The first stage has separate classifiers for each scale and aspect ratio; the second stage ranks all proposals from the previous stage. All SVMs are trained using structured output learning to score windows higher that overlap more with objects. Because the cascade is trained and tested over the same set of categories, it is unclear how well this approach generalises across categories.
- **RandomizedSeeds** [45] uses multiple randomised SEED superpixel maps [46] to score each candidate window. The scoring is done using a simple metric similar to “superpixel straddling” from `Object-`

ness, no additional cues are used. The authors show that using multiple superpixel maps significantly improves recall.

2.3 Alternative proposal methods

- **ShapeSharing** [47] is a non-parametric, data-driven method that transfers object shapes from exemplars into test images by matching edges. The resulting regions are subsequently merged and refined by solving graph cuts.
- **Multibox** [9], [48] trains a neural network to directly regress a fixed number of proposals in the image without sliding the network over the image. Each of the proposals has its own location bias to diversify the location of the proposals. The authors report top results on ImageNet.

2.4 Baseline proposal methods

We additionally consider a set of baselines that serve as reference points. Like all evaluated methods described earlier, the following baselines are class independent:

- **Uniform**[†]: To generate proposals, we uniformly sample the bounding box centre position, square root area, and log aspect ratio. We estimate the range of these parameters on the PASCAL VOC 2007 training set after discarding 0.5% of the smallest and largest values, so that our estimated distribution covers 99% of the data.
- **Gaussian**[†]: Likewise, we estimate a multivariate Gaussian distribution for the bounding box centre position, square root area, and log aspect ratio. After calculating mean and covariance on the training set we sample proposals from this distribution.
- **SlidingWindow**[†]: We place windows on a regular grid as is common for sliding window object detectors. The requested number of proposals is distributed across window sizes (width and height), and for each window size, we place the windows uniformly. This procedure is inspired by the implementation of Bing [18], [41].
- **Superpixels**[†]: As we will show, superpixels have an important influence on the behaviour of proposal methods. Since five of the evaluated methods build on [37], we use it as a baseline: each low-level segment is used as a detection proposal. This method serves as a lower-bound on recall for methods using superpixels.

It should be noted that with the exception of *Superpixels*, all the baselines generate proposal windows independent of the image content. *SlidingWindow* is deterministic given the image size (similar to *CrackingBing*), while the *Uniform* and *Gaussian* baselines are stochastic.

2.5 Proposals versus cascades

Many proposal methods utilise image features to generate candidate windows. One can interpret this process as a discriminative one; given such features a method quickly determines whether a window should be considered for

detection. Indeed, many of the surveyed methods include some form of discriminative learning (*SelectiveSearch* and *EdgeBoxes* are notable exceptions). As such, proposal methods are related to cascades [2], [49]–[51], which use a fast but inaccurate classifier to discard a vast majority of unpromising proposals. Although traditionally used for class specific detection, cascades can also apply to sets of categories [52], [53].

The key distinction between traditional cascades and proposal methods is that the latter is required to generalise beyond object classes observed during training. So what allows discriminatively trained proposal methods to generalise to unseen categories? A key assumption is that training a classifier for a large enough number of categories is sufficient to generalise to unseen categories (for example, after training on cats and dogs proposals may generalise to other animals). Additionally, the discriminative power of the classifier is often limited (e.g. Bing and Zhang), thus preventing overfitting to the training classes and forcing the classifier to learn coarse properties shared by all object (e.g. “objects are roundish”). This key distinction is also noted in [54]. We test the generalisation of proposal methods by evaluating on datasets with many additional classes in §4.

2.6 Controlling the number of proposals

In this work we will perform an extensive apples-to-apples comparison of the 12 methods (plus 4 baselines) listed in table 1. In order to be able to compare amongst methods, for each method we need to control the number of proposals produced per image. By default, the evaluated methods provide variable numbers of detection proposals, ranging from just a few ($\sim 10^2$) to a large number ($\sim 10^5$). Additionally, some methods output sorted or scored proposals, while others do not. Having more proposals increases the chance for high recall, thus for each method in all experiments we attempt to carefully control the number of generated proposals. Details are provided next.

Albeit not all having explicit control over the number of proposals, *Objectness*, *CPMC*, *Endres*, *SelectiveSearch*, *Rahtu*, *Bing*, *MCG*, and *EdgeBoxes* do provide scored or sorted proposals so we can use the top k . *Rantalankila*, *Rigor*, and *Geodesic* provide neither direct control over the number of proposals nor sorted proposals, but indirect control over k can be obtained by altering other parameters. Thus, we record the number of produced proposals on a subset of the images for different parameters and linearly interpolate between the parameter settings to control k . For *RandomizedPrim's*, which lacks any control over the number of proposals, we randomly sample k proposals.

Finally, we observed a number of methods produce duplicate proposals. All such duplicates were removed.

3 PROPOSAL REPEATABILITY

Training a detector on detection proposals rather than on all sliding windows modifies the appearance distribution of both positive and negative windows. In section 4, we look into how well the different object proposals overlap with ground truth annotations of objects, which is an analysis of the positive window distribution. In this section

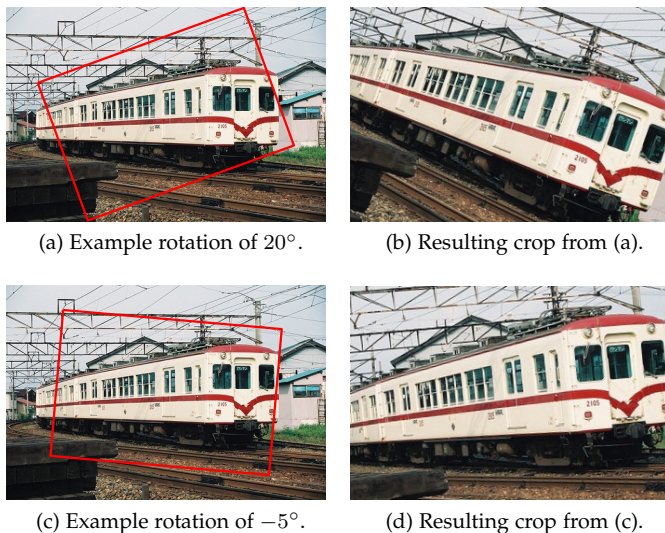


Figure 2: Examples of rotation perturbation. (a) shows the largest rectangle with the same aspect as the original image that can fit into the image under a 20° rotation, and (b) the resulting crop. All other rotations are cropped to the same dimensions, e.g. the -5° rotation in (c) to the crop in (d).

we analyse the distribution of negative windows: if the proposal method does not consistently propose windows on similar image content without objects or with partial objects, the classifier may have difficulty generating scores on negative windows on the test set. As an extreme, motivational example, consider a proposal method that generates proposals containing only objects on the training set but containing both objects and negative windows on the test set. A classifier trained on such proposals would be unable to differentiate objects from background, thus at test time would give useless scores for the negative windows. Thus we expect that a consistent appearance distribution for proposals *on the background* is likewise relevant for a detector.

We call the property of proposals being placed on similar image content the *repeatability* of a proposal method. Intuitively proposals should be repeatable on slightly different images with the same content. To evaluate repeatability we compare proposals that are generated for one image with proposals generated for a slightly modified version of the same image. PASCAL VOC [4] does not contain suitable images. An alternative is the dataset of [31], but it only consists of 54 images and even fewer objects. Instead, we opt to apply synthetic transformations to PASCAL images.

3.1 Evaluation protocol for repeatability

Our evaluation protocol is inspired by [31], which evaluates interest point repeatability. For each image in the PASCAL VOC 2007 test set [4], we generate several perturbed versions. We consider blur, rotation, scale, illumination, JPEG compression, and “salt and pepper” noise (see figures 3-4).

For each pair of reference and perturbed images we compute detection proposals with a given method (generating 1000 windows per image). The proposals are projected back from the perturbed into the reference image and then

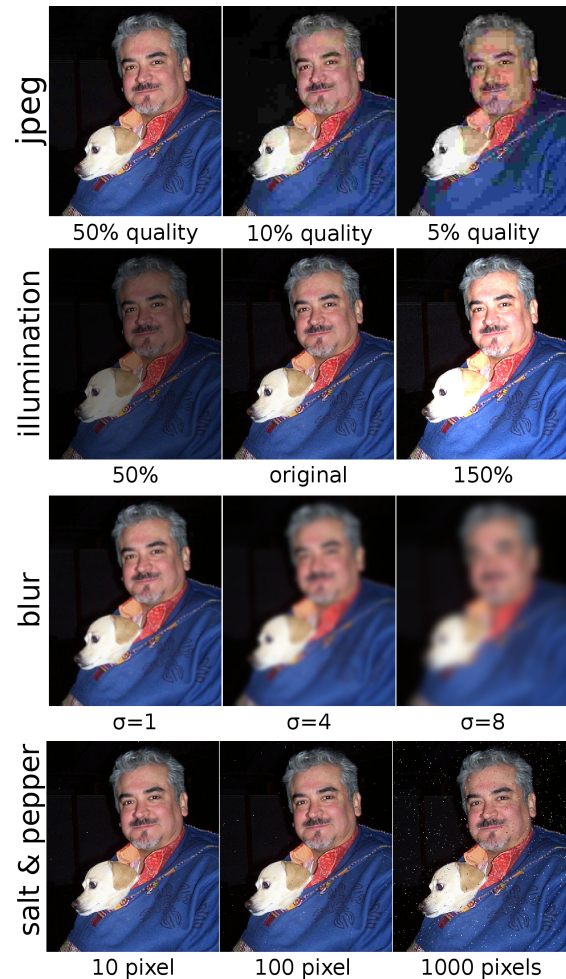


Figure 3: Illustration of the perturbation ranges used for the repeatability experiments.

matched to the proposals in the reference image. In the case of rotation, all proposals whose centre lies outside the image after projection are removed before matching. For matching we use the intersection over union (IoU) criterion and we solve the resulting bipartite matching problem greedily for efficiency reasons. Given the matching, we plot the recall for every IoU threshold and *define the repeatability to be the area under this “recall versus IoU threshold” curve between IoU 0 and 1*³. This is similar to computing the average best overlap (ABO, see §A) for the proposals on the reference image. Methods that propose windows at similar locations at high IoU—and thus on similar image content—are more repeatable, since the area under the curve is larger.

One issue regarding such proposal matching is that large windows are more likely to match than smaller ones since the same perturbation will have a larger relative effect on smaller windows. This effect is important to consider since different methods have very different distributions of proposal window sizes as can be seen in figure 5a. To reduce the impact of this effect, we bin the original image windows by area into 10 groups, and evaluate the area under the

³ In contrast to the average recall (AR) used in later sections, we use the area under the entire curve. We are interested in how much proposals change, which is independent of the PASCAL overlap criterion.

recall versus IoU curve per size group. In figure 5b we show the recall versus IoU curve for a small blur perturbation for each of the 10 groups. As expected, large proposals have higher repeatability. In order to measure repeatability independently of the distribution of window sizes, in all remaining repeatability experiments in figure 5 we show the (unweighted) average across the 10 size groups.

We omit the slowest two methods, *CPMC* and *Endres*, due to computational constraints (these experiments require running the detectors ~ 50 times on the entire PASCAL test set, once for every perturbation).

3.2 Repeatability experiments and results

There are some salient aspects of the result curves in figure 5 that need additional explanation. First, not all methods have 100% repeatability when there is no perturbation. This is due to random components in the selection of proposals for several methods. Attempting to remove a method’s random component is beyond the scope of this work and could potentially considerably alter the method. A second important aspect is the large drop of repeatability for most methods, even for subtle image changes. We observed that many of the methods based on superpixels are particularly prone to such perturbations. Indeed the *Superpixels* baseline itself shows high sensitivity to perturbations, so the instability of the superpixels likely explains much of this effect. Inversely we notice that methods that are not based on superpixels are most robust to small image changes (e.g. *Bing* and also the baselines that ignore image content).

We now discuss the details and effects of each perturbation on repeatability, shown in figure 5:

Scale (5c): We uniformly sample the scale factor from $.5\times$ to $2\times$, and test additional scales near the original resolution ($.9\times$, $.95\times$, $.99\times$, $1.01\times$, $1.05\times$, $1.1\times$). Upscaling is done with bicubic interpolation and downscaling with anti-aliasing. All methods except *Bing* show a drastic drop with small scale changes, but suffer only minor degradation for larger changes. *Bing* is more robust to small scale changes; however, it is more sensitive to larger changes due to its use of a coarse set of box sizes while searching for candidates (this also accounts for its dip in repeatability at half scales). The *SlidingWindow* baseline suffers from the same effect.

JPEG artefacts (5d): To create JPEG artefacts we write the target image to disk with the Matlab function `imwrite` and specify a quality settings ranging from 5% to 100%, see figure 3. Even the 100% quality setting is lossy, so we also include a lossless setting for comparison. Similar to scale change, even slight compression has a large effect and more aggressive compression shows monotonic degradation. Despite using gradient information, *Bing* is most robust to these kind of changes.

Rotation (5e): We rotate the image in 5° steps between -20° and 20° . To ensure roughly the same content is visible under all rotations, we construct the largest box with the same aspect as the original image that fits into the image under a 20° rotation and use this crop for all other rotations, see figure 2. All proposal methods are equally affected by image rotation. The drop of the *Uniform* and *Gaussian* baselines indicate the repeatability loss due to the fact that we are matching rotated bounding boxes.

Illumination (5f): To synthetically alter illumination of an image we changed its brightness channel in HSB colour space. We vary the brightness between 50% and 150% of the original image so that some over and under saturation occurs, see figure 3. Repeatability under illumination changes shows a similar trend as under JPEG artefacts. Methods based on superpixels are heavily affected. *Bing* is more robust, likely due to use of gradient information which is known to be fairly robust to illumination changes.

Blur (5g): We blur the images with a Gaussian kernel with standard deviations $0 \leq \sigma \leq 8$, see figure 3. The repeatability results again exhibit a similar trend although the drop is stronger for a small σ .

Salt and pepper noise (5h): We sample between 1 and 1000 random locations in the image and change the colour of the pixel to white if it is a dark and to black otherwise, see figure 3. Surprisingly, most methods already lose some repeatability when even a single pixel is changed. Significant degradation in repeatability for the majority of the methods occurs when merely ten pixels are modified.

Discussion: Small changes to an image cause noticeable differences in the set of detection proposals for all methods except *Bing*. The higher repeatability of *Bing* is explained by its sliding window pattern, which has been designed to cover almost all possible annotations with $\text{IoU} = 0.5$ (see also *Cracking Bing* [41]). As one cause for poor repeatability we identify the segmentation algorithm on which many methods build. Among all proposal methods, *EdgeBoxes* also performs favourably, possibly because it avoids the hard decision of grouping pixels into superpixels.

We also experimented with repeatability of boxes that touch annotations sufficiently ($\text{IoU} \geq 0.5$), which showed very similar trends, indicating that the issue of repeatability also applies to proposals that partially cover objects.

Different applications will be more or less sensitive to repeatability. Our results indicate that if repeatability is a concern, the proposal method should be selected with care. For object detection, another aspect of interest is recall, which we explore in the next section.

4 PROPOSAL RECALL

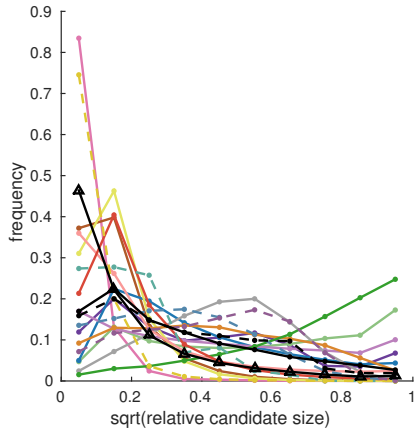
When using detection proposals for detection it is important to have a good coverage of the objects of interest in the test image, since missed objects cannot be recovered in the subsequent classification stage. Thus it is common practice to evaluate the quality of proposals based on the recall of the ground truth annotations.

4.1 Evaluation protocol for recall

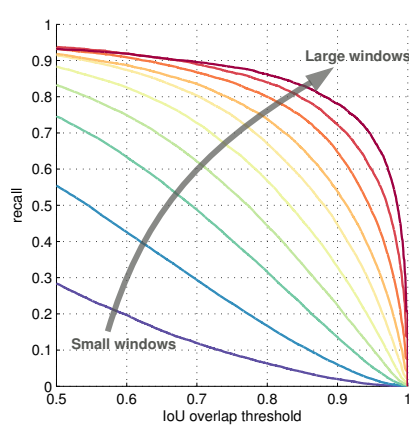
The protocol introduced in [12] (using the PASCAL VOC 2007 dataset [4]) has served as a guideline for most evaluations in the literature. While previous papers do show various comparisons on PASCAL, the train and test sets vary amongst papers, and the metrics shown tend to favour different methods. We provide an extensive and unified evaluation and show that different metrics result in different rankings of proposal methods (e.g. see figure 6b versus 7b).



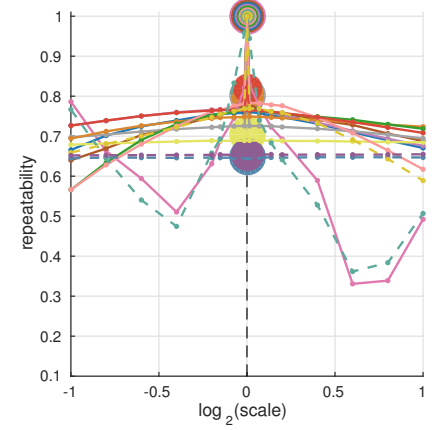
Figure 4: Example of the image perturbations considered. Top to bottom, left to right: original, blur, illumination, JPEG artefact, rotation, scale perturbations, and “salt and pepper” noise.



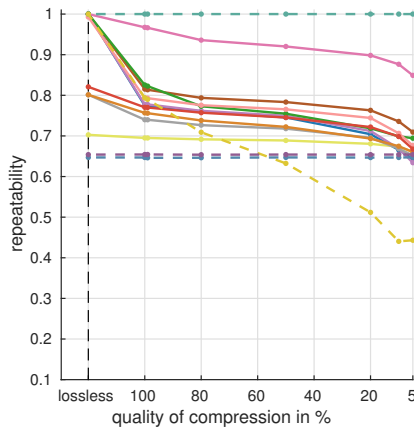
(a) Histogram of proposal sizes on PASCAL.



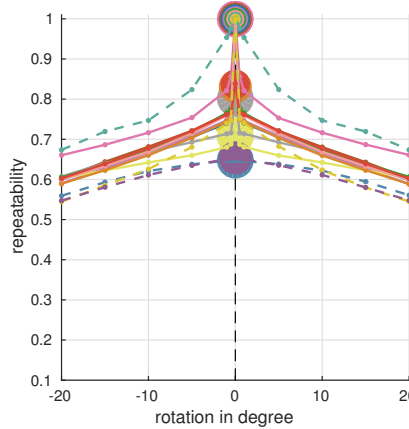
(b) Example of recall at different scales.



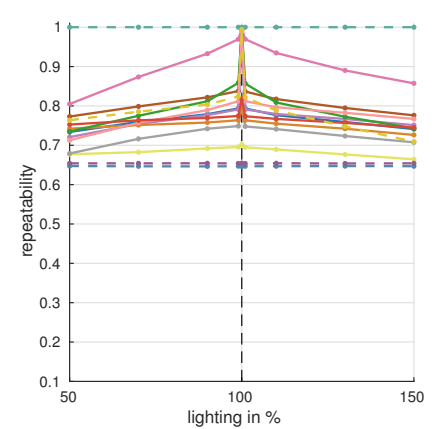
(c) Scale.



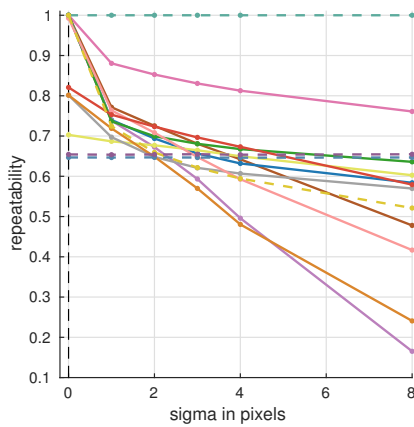
(d) JPEG artefacts.



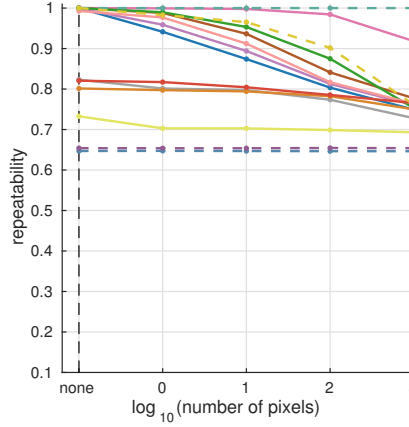
(e) Rotation.



(f) Illumination.



(g) Blur.



(h) Salt and pepper noise.

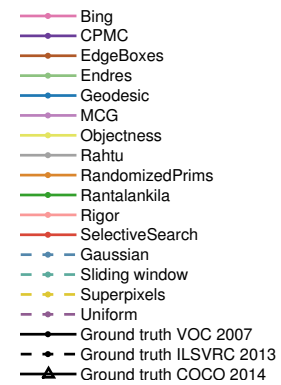


Figure 5: Repeatability results under various perturbations.

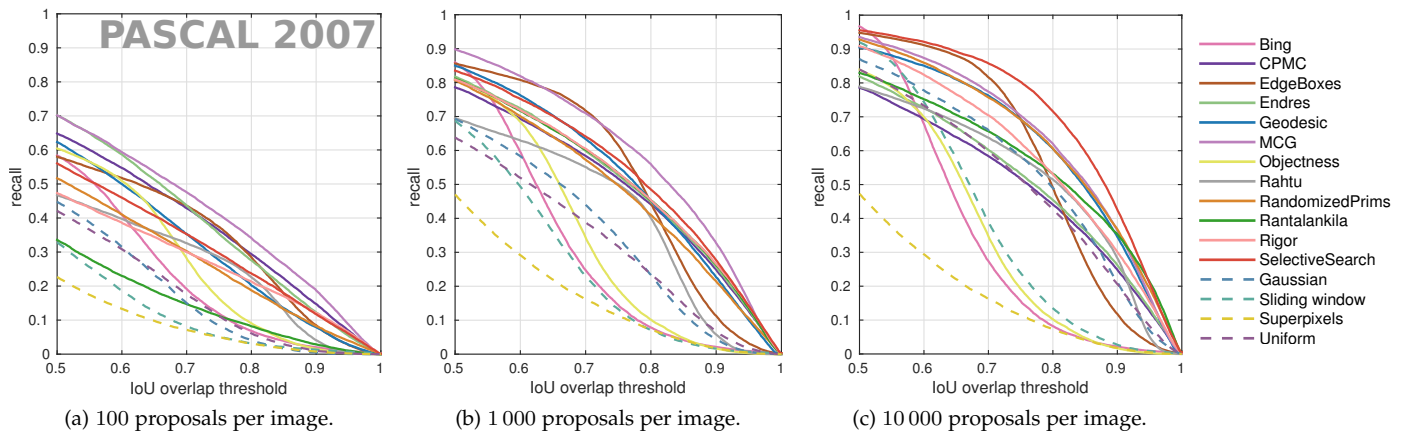


Figure 6: Recall versus IoU threshold on the PASCAL VOC 2007 test set.

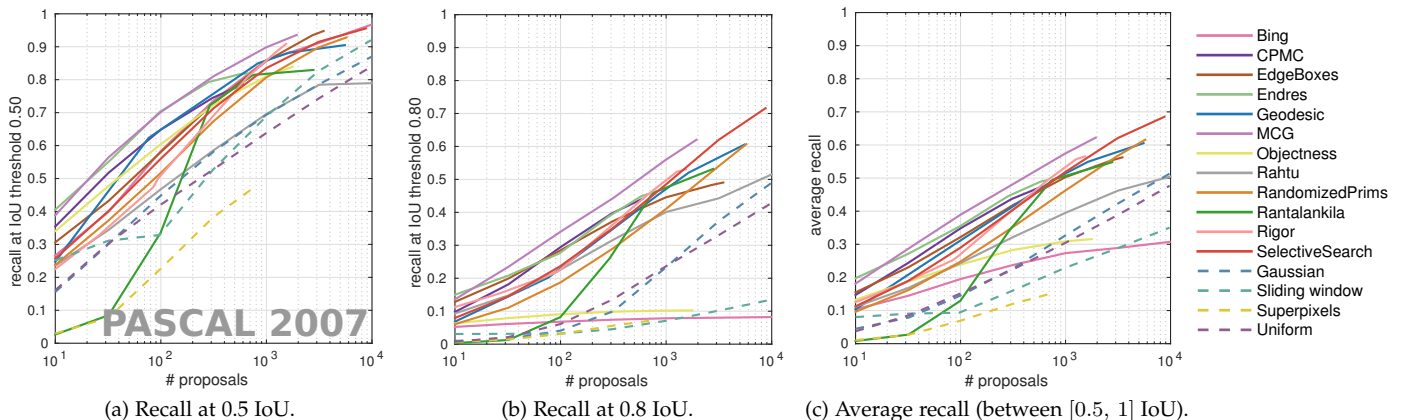


Figure 7: Recall versus number of proposal windows on the PASCAL VOC 2007 test set.

Metrics: Evaluating (class agnostic) detection proposals is quite different from traditional class-specific detection [55] since most metrics (class confusion, background confusion, precision, etc.) do not apply. Instead, one of the primary metrics for evaluating proposals is, for a fixed number of proposals, the fraction of ground truth annotations covered as the intersection over union (IoU) threshold is varied (figure 6). Another common and complementary metric is, for a fixed IoU threshold, proposal recall as the number of proposals is varied (figure 7a, 7b). Finally, we define and report a novel metric, the average recall (AR) between IoU 0.5 to 1, and plot AR versus number of proposals (figure 7c).

PASCAL: We evaluate recall on the full PASCAL VOC 2007 test set [4], which includes 20 object categories present in ~ 5000 unconstrained images. For the purpose of proposal evaluation we include all 20 object categories and all ground truth bounding boxes, including “difficult” ones, since our goal is to measure maximum recall. In contrast to [12], we compute a matching between proposals and ground truth, so one proposal cannot cover two objects. Note that while different methods may be trained on different sets of object categories and subsets of data, we believe evaluating on all categories at test time is appropriate as we care about absolute proposal quality. Such an evaluation strategy is further supported as many methods have no training stage, yet provide competitive results (e.g. *SelectiveSearch*).

ImageNet: The PASCAL VOC 2007 test set, on which most proposal methods have been previously evaluated, has only 20 categories, yet detection proposal methods claim to predict proposals for *any* object category. Thus there is some concern that the proposal methods may be tuned to the PASCAL categories and not generalise well to novel categories. To investigate this potential bias, we also evaluate methods on the larger ImageNet [5] 2013 validation set, which contains annotations for 200 categories in over ~ 20000 images. It should be noted that these 200 categories are *not* fine grained versions of the PASCAL ones. They include additional types of animals (e.g. crustaceans), food items (e.g. hot-dogs), household items (e.g. diapers), and other diverse object categories.

MS COCO: Although ImageNet has 180 more classes than PASCAL, it is still similar in statistics like number of objects per image and size of objects. Microsoft Common Objects in Context (MS COCO) [6] has more objects per image, smaller objects, but also fewer object classes (80 object categories). We evaluate the recall of this dataset to further investigate potential biases of proposal methods. We evaluate the recall on all annotations excluding the “crowd” annotations which may mark large image areas including a lot of background.

4.2 Recall results

PASCAL Results in figure 6 and 7 present a consistent trend across the different metrics. *MCG*, *EdgeBoxes*, *Selective-*

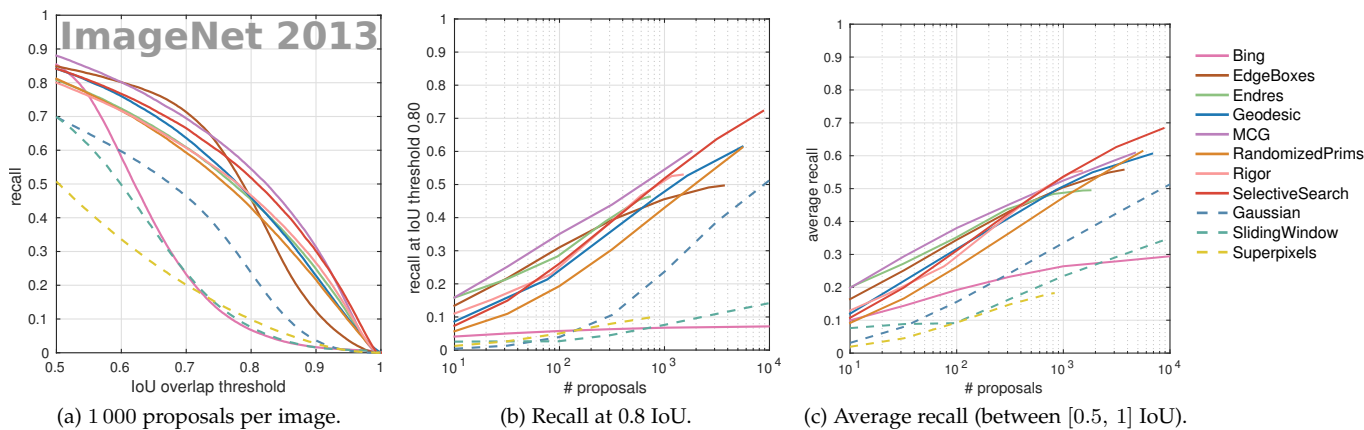


Figure 8: Recall on the ImageNet 2013 validation set.

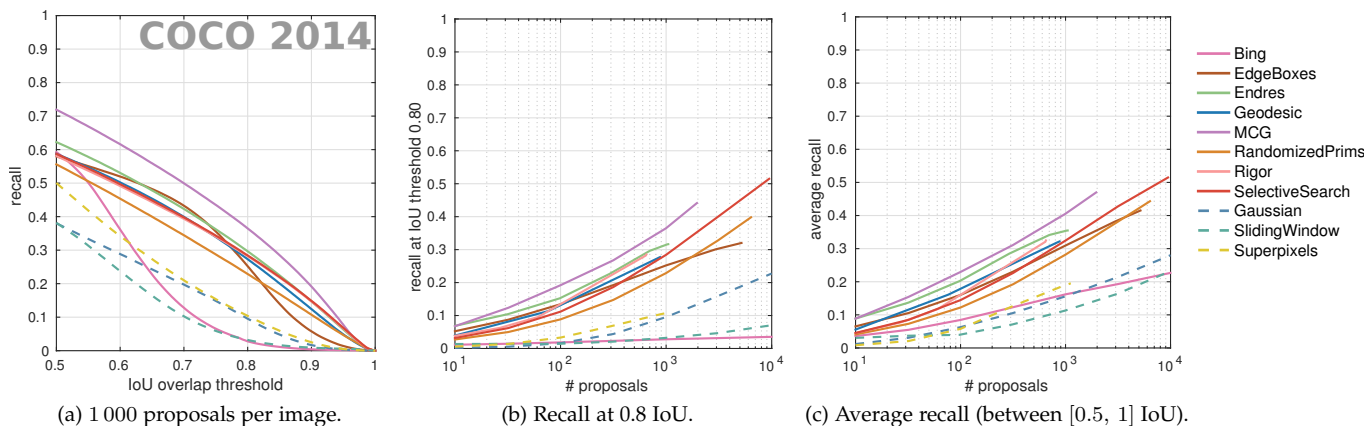


Figure 9: Recall on the MS COCO 2014 validation set.

Search, Rigor, and Geodesic are the best methods across different numbers of proposals. SelectiveSearch is surprisingly effective despite being fully hand-crafted (no machine learning involved). When considering less than 10^3 proposals, MCG, Endres, and CPMC provide strong results.

Overall, the methods fall into two groups: well localised methods that gradually lose recall as the IoU threshold increases and methods that only provide coarse bounding box locations, so their recall drops rapidly. All baseline methods, as well as Bing, Rahtu, Objectness, and EdgeBoxes fall into the latter category. Bing in particular, while providing high repeatability, only provides high recall at IoU = 0.5 and drops dramatically when requiring higher overlap (the reason for this is identified in [41]).

Baselines: When inspecting figure 6 from left to right, one notices that with few proposals the baselines provide relatively low recall (figure 6a). However as the number of proposals increases, Gaussian and Uniform become more competitive (figure 6b). In relative gain, detection proposal methods have most to offer for low numbers of windows.

Average Recall: Rather than reporting recall at particular IoU thresholds, we also report the average recall (AR) between IoU 0.5 to 1 (which is related to the ABO metric, see §A), and plot AR for varying number of proposals in figure 7c. Much like the average precision (AP) metric for

(class specific) object detection, AR summarises proposal performance across IoU thresholds (for a given number of proposals). In fact, in §5 we will show that AR correlates well with detection performance. As can be seen in figure 7c, MCG performs well across the entire range of number of proposals. Endres and EdgeBoxes work well for a low number of proposals while for a higher number of proposals Rigor and SelectiveSearch perform best.

ImageNet: As discussed, compared to PASCAL, ImageNet includes $10\times$ ground truth classes and $4\times$ images. Somewhat surprisingly the ImageNet results in figure 8 are almost identical to the ones in figures 6b, 7b, and 7c. To understand this phenomenon, we note that the statistics of ImageNet match the ones of PASCAL. In particular the typical image size and the mean number of object annotation per image (three) is similar in both datasets. This helps explain why the recall behaviour is similar, and why methods tuned on PASCAL still perform well on ImageNet.

MS COCO: We present the same results for MS COCO in figure 9. We see different absolute numbers, yet similar trends with some notable exceptions as can be seen in figure 10a. EdgeBoxes no longer ranks significantly better than SelectiveSearch, Geodesic and Rigor for few proposals. MCG and Endres improve relative to the other methods, in particular for a higher number of proposals. We

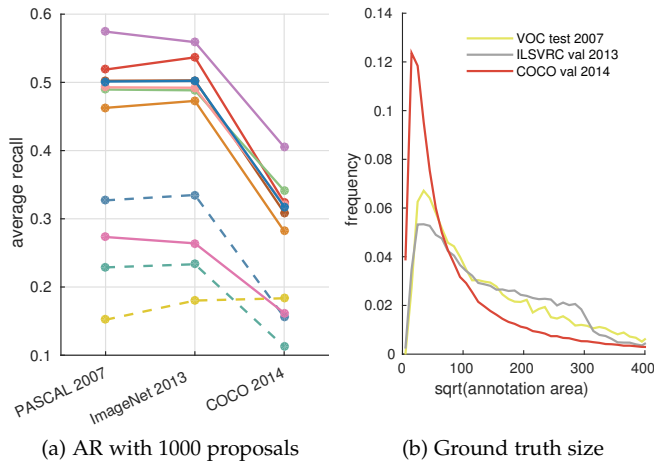


Figure 10: Comparison between all considered datasets: PASCAL VOC 2007 test set, ImageNet 2013 validation set, MS COCO 2014 validation set (see methods legend fig. 7c).

attribute these difference to different statistics of the dataset, particularly the different size distribution, see figure 10b.

Overall, MCG is the top performing method across all datasets in terms of both recall and AR at all settings. This is readily apparent in figure 10a.

Generalisation: We emphasise that although the results on PASCAL, ImageNet, and MS COCO are quite similar (see figure 10a), ImageNet covers 200 object categories, many of them unrelated to the 20 PASCAL categories and COCO has significantly different statistics. In other words, there is no measurable over-fitting of the detection proposal methods towards the PASCAL categories. This suggests that proposal methods transfer adequately amongst object classes, and can thus be considered true “objectness” measures.

5 USING THE DETECTION PROPOSALS

In this section we analyse detection proposals for use with object detectors. We consider two well known and quite distinct approaches to object detection. First we use a variant of the popular DPM part-based sliding window detector [3], specifically the LM-LLDA detector [56]. We also test the state of the art R-CNN [8] and Fast R-CNN [16] detectors which couple object proposals with a convolutional neural network classification stage. Our goals are twofold. First, we aim to measure the performance of different proposal methods for object detection. Second, we are interested in evaluating how well the proposal metrics reported in the previous sections can serve as a proxy for predicting final detection performance. All following experiments involving proposals use 1 000 proposals.

5.1 Detector responses around objects

As a preliminary experiment, we aim to quantify the importance of having well localised proposals for object detection. We begin by measuring how detection scores are affected by the overlap between the detector window and the ground truth annotation on the PASCAL 2007 test set [4]. When considering the detectors’ bounding box prediction, we use the refined position to compute the overlap.

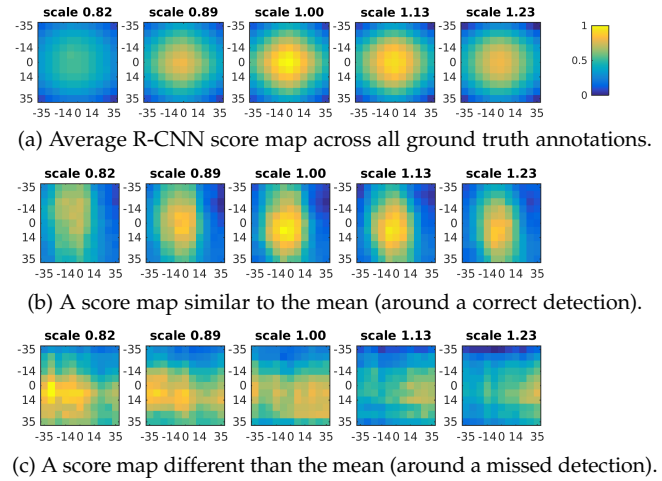


Figure 11: Normalised score maps of the R-CNN around ground truth annotations on the PASCAL 2007 test set. One grid cell in each map has width and height of ~ 7 px after the object height has been resized to the detector window of 227×227 px (3% of the object height).

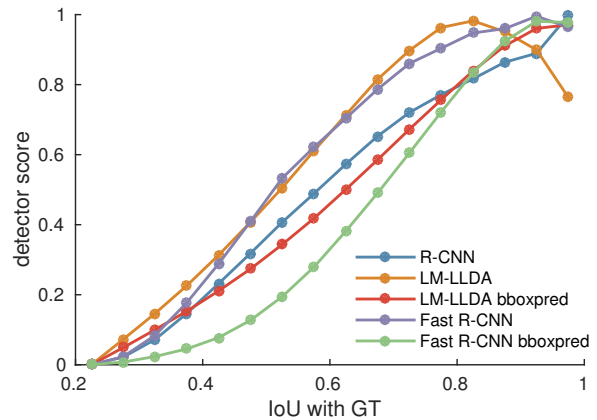


Figure 12: Normalised detector scores as a function of the overlap between the detector window and the ground truth.

Score map: Figure 11a shows the average R-CNN detection score around the ground truth annotations. We notice that the score map is symmetric and attains a maximum at the ground truth object location. In other words, the detector has no systematic spatial or scale bias. However, averaging the score maps removes small details and imperfections of individual score maps. When considering individual activations instead of the average, we observe a high variance in the quality of the score maps, see figures 11b and 11c.

Score vs IoU: In figure 12 we show average detection scores for proposals with varying IoU overlap with the ground truth. The scores have been scaled between zero and one per class before averaging across classes. The drop of the LM-LLDA scores at high overlaps is due to a bias introduced during training by the latent location estimation on positive samples; this bias is compensated for by the subsequent bounding box prediction stage of LM-LLDA. For Fast R-CNN, the bounding box prediction effectively improves proposal IoU with the ground truth and results in a substantial shift of the curve to the right.

	aero	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
LM-LLDA Dense	33.7	61.3	12.4	18.5	26.7	53.0	57.2	22.4	22.7	25.6	25.1	14.0	59.2	51.0	39.1	13.6	21.7	38.0	48.8	44.0	34.4
Bing	-7.5	-23.2	-6.2	-8.1	-10.6	-13.3	-17.5	-6.8	-9.8	-15.4	-7.5	-1.4	-19.6	-19.0	-16.1	-3.4	-6.6	-18.1	-18.8	-10.0	-11.9
CPMC	-1.0	-15.0	-0.2	-4.4	-13.5	-1.8	-9.2	3.2	-9.1	-2.6	5.1	2.2	-4.2	-4.8	-7.0	-2.0	-2.6	1.2	-4.1	-4.9	-3.7
EdgeBoxes	-2.0	-6.1	-0.7	-3.8	-6.7	0.6	-5.8	-1.1	-2.0	-1.8	-4.6	0.4	-1.3	-1.3	-3.0	-1.7	-0.1	-0.9	-0.2	-1.1	-2.2
Endres	-1.5	-5.8	-0.6	-4.8	-12.7	-1.1	-7.1	3.4	-6.9	-3.2	4.7	1.9	-2.4	-2.4	-7.7	-2.8	-1.9	1.5	0.4	-4.2	-2.7
Geodesic	-1.9	-8.1	-0.2	-4.6	-14.4	0.6	-6.5	2.6	-7.3	-1.3	4.7	2.4	-2.5	-2.7	-4.7	-1.2	-0.7	-0.1	1.9	0.2	-2.2
MCG	-0.7	-7.2	0.1	-3.6	-6.7	-1.2	-7.0	3.4	-3.2	-2.3	5.0	1.9	-3.5	-1.3	-1.5	-1.1	-1.3	2.2	0.3	0.5	-1.4
Objectness	-10.3	-15.1	-2.0	-6.2	-11.0	-9.5	-13.0	-3.6	-10.0	-6.4	-7.8	-1.0	-11.6	-15.9	-13.0	-2.7	-5.8	-11.2	-10.9	-12.9	-9.0
Rahtu	-0.3	-13.2	-0.3	-1.2	-13.0	-0.6	-12.0	3.3	-10.5	-4.3	2.0	2.1	-3.2	-4.9	-7.9	-2.8	-4.9	-5.0	0.0	-3.7	-4.0
Rand.Prim	2.1	-10.4	-0.5	-4.5	-13.2	-1.9	-10.1	5.0	-6.7	-3.5	2.0	2.4	-4.4	-5.1	-10.0	-2.3	-1.8	1.2	-3.8	-4.4	-3.5
Rantalankila	0.5	-13.6	0.3	-3.0	-12.9	-3.6	-9.0	4.4	-5.6	-3.7	4.1	2.5	-2.2	-4.0	-7.8	-2.5	-3.8	2.1	-1.5	-0.7	-3.0
Rigor	1.7	-7.9	0.5	-4.1	-12.4	-0.8	-9.0	6.3	-6.9	-1.7	1.8	2.9	-0.9	-3.3	-7.7	-1.8	-1.3	1.6	-1.2	-1.7	-2.3
SelectiveSearch	1.3	-7.7	1.0	-4.3	-11.1	-1.7	-7.8	3.9	-4.8	-1.5	5.4	2.2	-1.4	-3.8	-6.0	-1.5	-0.8	0.6	-2.4	-2.1	-2.1
Gaussian	-6.6	-13.4	-0.7	-4.4	-15.0	-6.1	-16.0	0.9	-9.1	-8.0	0.3	1.2	-4.2	-6.9	-10.3	-2.3	-6.5	-4.5	-3.6	-12.1	-6.4
SlidingWindow	-21.8	-20.7	-3.2	-8.1	-16.6	-14.7	-22.1	-0.7	-9.8	-11.7	-10.2	-1.4	-14.7	-20.1	-14.8	-3.8	-7.7	-21.0	-20.8	-14.8	-12.9
Superpixels	-23.9	-52.2	-3.1	-9.4	-17.4	-43.9	-42.3	-10.2	-11.3	-12.6	-15.8	-8.5	-50.1	-41.7	-30.9	-4.4	-10.6	-25.2	-39.7	-8.2	-23.1
Uniform	-3.2	-18.8	-4.0	-4.8	-15.2	-8.6	-16.6	0.2	-10.4	-8.8	3.7	1.3	-6.6	-11.3	-10.2	-3.6	-8.9	-5.8	-5.1	-20.2	-7.8
Top methods avg.	-0.3	-7.4	0.1	-4.1	-10.2	-0.5	-7.2	3.0	-4.8	-1.7	2.5	2.0	-1.9	-2.5	-4.6	-1.5	-0.8	0.7	-0.3	-0.8	-2.0

Table 2: LM-LLDA detection results on PASCAL 2007 (with bounding box regression). The top row indicates the average precision (AP) of LM-LLDA alone, while the other rows show the difference in AP when adding proposal methods. Green indicates improvement of at least 2 AP, blue indicates minor change ($-2 \leq \Delta AP < 2$), and white indicates a decrease by more than 2 AP. EdgeBoxes achieves top results on 6 of the 20 categories; MCG performs best overall with -1.4 mAP loss.

Proposals	LM-LLDA	R-CNN	Fast R-CNN	Δ Train
Dense	33.5/34.4	-	-	-
Bing	21.8/22.4	36.7	37.3/49.0	+6.3
CPMC	30.0/30.7	51.7	53.7/57.1	-1.3
EdgeBoxes	31.8/32.2	53.0	55.4/60.4	+3.3
Endres	31.2/31.7	52.8	54.2/57.4	-0.2
Geodesic	31.8/32.2	53.8	53.6/57.5	-0.4
MCG	32.5/33.0	56.5	58.1/60.3	+1.8
Objectness	25.0/25.4	39.7	41.5/51.4	+9.1
Rahtu	29.6/30.4	46.1	48.9/53.6	+0.7
RandomizedPrims	30.5/30.9	51.6	53.2/57.6	-0.6
Rantalankila	30.9/31.4	53.1	55.0/57.9	-0.5
Rigor	31.5/32.1	54.1	55.4/58.4	-0.2
SelectiveSearch	31.7/32.3	54.6	56.3/59.5	+0.0
Gaussian	27.3/28.0	40.6	44.6/50.8	+0.8
Sliding window	20.7/21.5	32.7	32.7/44.8	+3.3
Superpixels	11.2/11.3	17.6	15.4/20.3	-2.0
Uniform	26.0/26.6	37.3	39.5/46.9	-0.1

Table 3: Mean average precision (mAP) on PASCAL 2007 for multiple detectors and proposal methods (using 1000 proposals). LM-LLDA and Fast R-CNN results shown before/after bounding box regression. The final column shows the change in mAP obtained from re-training Fast R-CNN (with box regression) for the specific proposal method.

Localisation: We observe from figure 12 that both LM-LLDA and R-CNN exhibit an almost linear increase in detection score as IoU increases (especially between 0.4 and 0.8 IoU). From this we conclude that there is no IoU threshold that is “sufficiently good” for obtaining top detection quality. We thus consider that improving localisation of proposals is as important as increasing ground truth recall, and the linear relation helps motivate us to linearly reward localisation in the average recall metric (see §4.2). For Fast R-CNN there is also an almost linear relation, but performance saturates earlier. Thus, Fast R-CNN is likely to also benefit from better localisation, but up to a point.

5.2 LM-LLDA detection performance

We use pre-trained LM-LLDA [56] models to generate dense detections using the standard sliding window setup and subsequently apply different proposals to filter these

detections at test time. This does not speed-up detection, but enables evaluating the effect of proposals on detection quality. A priori we may expect that detection results will deteriorate due to lost recall, but conversely, they may improve if the proposals filter out windows that would otherwise be false positives.

Implementation: We take the raw detections of LM-LLDA before non-maximum suppression (nms) and filter them with the detection proposals of each method. We keep all detections that overlap more than 0.8 IoU with a candidate proposal and subsequently apply nms to the surviving detections. As a final step we do bounding box regression, as is common for DPM models [3]. Note that this procedure returns predictions near to, but distinct from, each proposal.

Results: Table 3, LM-LLDA columns, show that using 1000 proposals decreases detection quality compared with the original sliding window setup⁴ by about 1-2 mAP for the best performing methods, see top row (Dense) versus the rows below. The five top performing methods all have mAP between 32.0 and 33.0 and are marked in green: MCG, SelectiveSearch, EdgeBoxes, Geodesic, and Rigor. Note that the difference between these methods and the Gaussian baseline is fairly small (33.0 versus 28.0 mAP).

When we compare these results with figure 7c at 1000 proposals, we see that methods are ranked similarly. Methods with high average recall (AR) also have high mAP, and methods with lower AR also have lower mAP. We analyse the correlation between AR and mAP more closely in §5.4.

From table 2 we see that the per-class performance can be grouped into three cases: classes on which the best proposals (1) clearly hurt performance (bicycle, boat, bottle, car, chair, horse, mbike, person), (2) improve performance (cat, table, dog), (3) do not show significant change (all remaining classes). In the case of (1) we observe both reduced recall and reduced precision in the detection curves, probably because bad localisation decreases the scores of strong detections.

4. Not to be confused with the SlidingWindow proposals baseline.

	aero	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
Bing	56.6	54.9	45.0	28.6	24.6	53.9	63.5	72.5	15.6	59.4	49.0	59.7	68.5	60.3	50.7	16.5	49.0	42.8	64.8	44.9	49.0
CPMC	65.2	61.8	58.2	37.2	17.9	71.0	67.3	76.7	22.9	61.2	64.6	70.1	77.0	69.2	54.8	18.5	52.6	63.4	71.7	61.5	57.1
EdgeBoxes	67.0	69.9	59.8	46.1	28.3	72.9	72.3	73.8	28.8	68.1	62.4	67.6	79.2	73.6	62.4	28.2	55.8	61.2	70.4	59.7	60.4
Endres	61.5	70.8	57.1	33.5	18.0	72.5	68.8	77.3	21.7	61.8	64.5	68.2	78.0	69.9	56.2	21.4	54.5	63.2	72.4	56.9	57.4
Geodesic	63.2	68.0	55.9	39.2	19.8	71.1	70.4	74.4	24.8	65.0	63.5	65.6	78.7	69.2	58.0	20.4	54.5	57.8	70.2	60.9	57.5
MCG	66.6	69.1	60.1	42.0	28.5	71.9	72.3	77.3	30.2	61.3	62.4	69.8	77.4	68.2	62.2	27.5	57.6	66.0	75.8	59.4	60.3
Objectness	62.4	61.5	51.0	32.0	19.3	65.8	64.3	69.5	18.0	55.4	51.4	60.1	74.1	64.7	50.9	17.3	41.9	50.9	67.8	49.0	51.4
Rahtu	62.8	60.9	53.3	35.1	15.3	72.6	60.5	75.1	15.4	56.9	61.6	66.3	76.3	65.2	51.2	14.1	44.6	58.1	72.0	54.3	53.6
RandomizedPrims	70.2	68.2	55.5	39.5	18.5	72.3	63.7	76.8	25.7	62.4	64.2	68.7	76.6	68.5	51.0	22.4	53.1	62.9	72.4	59.7	57.6
Rantalankila	64.7	66.1	57.2	37.8	19.7	74.2	67.5	78.2	23.0	63.6	63.4	70.3	78.6	69.8	55.9	21.4	50.8	64.3	74.1	58.3	57.9
Rigor	62.6	70.5	57.5	40.1	15.9	72.9	65.7	77.9	28.6	65.1	63.7	68.6	77.9	68.9	54.8	23.3	56.3	63.8	73.7	60.3	58.4
SelectiveSearch	70.3	66.9	61.5	42.2	21.7	68.3	68.7	76.3	27.5	65.9	67.0	69.8	75.5	68.9	57.9	24.6	53.6	63.7	76.0	62.4	59.5
Gaussian	53.9	66.1	46.6	24.6	10.0	66.6	52.2	77.1	20.6	48.7	64.1	65.5	75.6	64.2	47.0	14.2	38.1	58.2	70.5	53.0	50.8
SlidingWindow	42.0	57.7	40.1	23.7	9.3	60.8	47.8	72.8	12.5	42.1	44.7	63.7	72.8	62.5	44.5	8.5	34.3	47.7	62.3	46.6	44.8
Superpixels	29.7	5.5	19.8	10.4	9.0	7.4	24.4	42.0	15.1	39.9	6.6	30.3	10.7	13.7	12.8	8.9	40.7	18.1	4.9	55.6	20.3
Uniform	51.0	58.0	38.6	24.6	11.7	64.3	50.9	72.3	14.8	43.4	62.6	63.4	73.9	59.3	43.4	10.8	27.5	60.4	69.0	38.3	46.9
best per class	70.3	70.8	61.5	46.1	28.5	74.2	72.3	78.2	30.2	68.1	67.0	70.3	79.2	73.6	62.4	28.2	57.6	66.0	76.0	62.4	62.1

Table 4: Fast R-CNN (model M) detection results (AP) on PASCAL VOC 2007. Bold numbers indicate the best proposal method per class, green numbers are within 2 AP of the best result. The “best per class” row shows the best performance when choosing the optimal proposals per class, improving from 60.4 mAP (EdgeBoxes) to 62.1 mAP.

5.3 R-CNN detection performance

The highly successful and widely used R-CNN detector [8] couples detection proposals with a convolutional neural network classification stage. It was designed from the ground up to rely on proposals, making it a perfect candidate for our case study. We report results for both the original R-CNN detector and also the improved Fast R-CNN [16]. We focus primarily on Fast R-CNN due to its efficiency and higher detection accuracy.

Implementation: For each proposal method we re-train and test Fast R-CNN (using the medium model M for efficiency). Unlike Fast R-CNN, the original R-CNN is fairly slow to train; we therefore experiment with the R-CNN model that is published with the code and which has been trained on 2 000 SelectiveSearch proposals.

Results: Although the absolute mAP numbers are considerably higher for Fast R-CNN (nearly double mAP), the results (Fast R-CNN and R-CNN) in table 3 show a similar trend than the LM-LLDA results. As expected, SelectiveSearch, with which Fast R-CNN was developed, performs well, but multiple other proposal methods get similar results. The five top performing methods are similar to the top methods for LM-LLDA: Rantalankila edges out EdgeBoxes for R-CNN and Geodesic for Fast R-CNN. EdgeBoxes and MCG provide the best results. The gap between Gaussian and the top result is more pronounced (60.4 versus 50.8 mAP), but this baseline still performs surprisingly well considering it disregards the image content. We show per-class Fast R-CNN results in table 4.

Retraining: To provide a fair comparison amongst proposal methods, the “Fast R-CNN” column in table 3 reports results after re-training for each method. The rightmost column of table 3 shows the change in mAP when comparing Fast R-CNN (with bounding box regression) trained with 1 000 SelectiveSearch proposals and applied at test time with a given proposal method, versus Fast R-CNN trained for the test time proposal method.

Most methods improve from re-training, although the performance of a few degrades. While in most cases the

change in mAP is within 1-2 points, re-training provided substantial benefits for Bing, EdgeBoxes, Objectness, and SlidingWindow. These methods all have poor localisation at high IoU (see figure 6); re-training likely allows Fast R-CNN to compensate for their inferior localisation.

Summary: We emphasise that the various proposal methods exhibit similar ordering with all tested detectors (LM-LLDA, R-CNN, and Fast R-CNN). Our experiments did not reveal any proposal methods as being particularly well-adapted for certain detectors; rather, for object detection some proposals methods are strictly better than others.

5.4 Predicting detection performance

We aim to determine which recall metrics from section 4 (figures 6 and 7) serve as the best predictor for detector performance. In figure 13 we show the Pearson correlation coefficient between detector performance and two recall metrics: recall at different overlap thresholds (left columns) and the *average recall* (AR) between IoU of 0.5 to 1.0 (right columns)⁵. As before, we use 1 000 proposals per method.

We begin by examining correlation between detection performance and recall at various IoU thresholds (figure 13, left columns). All detectors show a strong correlation (> 0.9) at an IoU range of roughly 0.6 to 0.8, with the exception of Fast R-CNN with bounding box prediction, which correlates better for lower overlap. Note that recall at IoU of 0.5 is actually only weakly correlated with detection performance, and methods that optimise for IoU of 0.5, such as Bing, are not well suited for use with object detectors (see table 3). Thus, although recall at IoU of 0.5 has been traditionally used to evaluate object proposals, our analysis shows that it is *not* a good metric for predicting detection performance.

The correlation between detection performance and AR is quite strong, see figure 13, right columns. Computing the AR over a partial IoU range (e.g. 0.6 to 0.8) can further increase the correlation; however, since the effect is generally

5. We compute the average between 0.5 and 1 IoU (and not between 0 and 1 as in §3), because we are interested in recall above the PASCAL evaluation criterion of 0.5 IoU. Proposals with worse overlap than 0.5 are not only harder to classify correctly, but require a potentially large subsequent location refinement to become a successful detection.

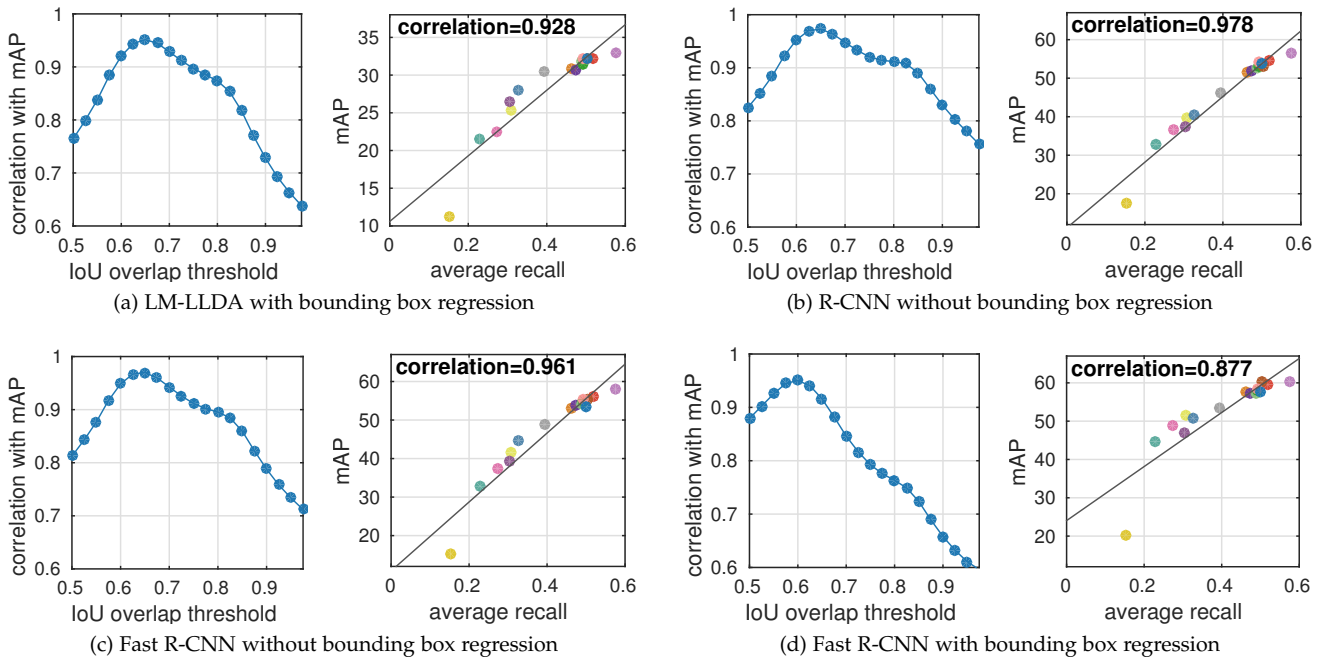


Figure 13: Correlation between detector performance on PASCAL 07 and different proposal metrics. Left columns: correlation between mAP and recall at different IoU thresholds. Right columns: correlation between mAP and AR.

minor, we opted to use AR over the entire range from 0.5 to 1.0 for simplicity. While the strong correlation does not imply that the AR can perfectly predict detection performance, as figure 13 shows, the relationship is surprisingly linear. AR over the full range of 0 to 1 IoU (which is similar to ABO, see appendix A) has weaker correlation with mAP, since proposals with low overlap are not sufficient for a successful detection under the PASCAL criterion and are also harder to classify.

For detectors with bounding box regression, the AR computation can be restricted to a tighter IoU range. In figure 12, we can observe that detection score of Fast R-CNN saturates earlier. Thus there is little benefit in proposals that are perfectly localised as the bounding box refinement improves the localisation of those proposals. If we restrict the AR to IoU from 0.5 to 0.7, we obtain a higher correlation of 0.949 for Fast R-CNN with bounding box regression (compared to 0.877 in figure 13d).

For a more detailed analysis of the correlation between mAP and AR we show the correlation for each class for different detectors in figure 14. The per-class correlation is highest for R-CNN and Fast R-CNN without regression.

We conclude that AR allows us to identify good proposal methods for object detection. The AR metric is simple, easy to justify, and is strongly correlated with detection performance. Note that our analysis only covers the case in which all methods produce the same number of proposals. As Girshick [16] points out, as the number of proposals increases, AR will necessarily increase but resulting detector performance saturates and may even degrade. For a fixed number of proposals, however, AR is a good predictor of detection performance. We suggest that future proposal methods should aim to optimise this metric.

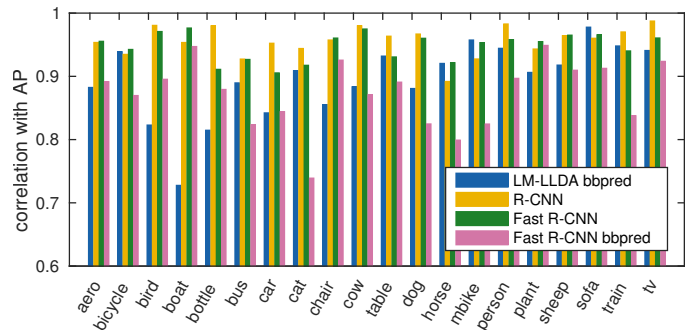


Figure 14: Correlation between AR and AP for each PASCAL VOC class and detector across all proposal methods.

5.5 Tuning proposal methods

All previous experiments evaluate proposal methods using original parameter settings. However many methods have free parameters that allow for fine-tuning. For example, when adjusting window sampling density and the non-maximum suppression (nms) in *EdgeBoxes* [20], it is possible to trade-off low recall with good localisation for higher recall with worse localisation (a similar observation was made in [40]). Figure 15 compares different versions of *EdgeBoxes* tuned to maximise recall at different IoU points Δ (we set $\alpha = \max(0.65, \Delta - 0.05)$, $\beta = \Delta + 0.05$, see [20] for details). *EdgeBoxes* tuned for $\Delta = 0.70$ or 0.75 maximises AR and also results in the best detection results.

While originally *EdgeBoxes* allowed for optimising recall for a particular IoU threshold, we consider a new variant that directly maximises AR (marked ‘AR’ in figure 15) to further explore the link between AR and detection quality. To do so, we alter its greedy nms procedure to make it

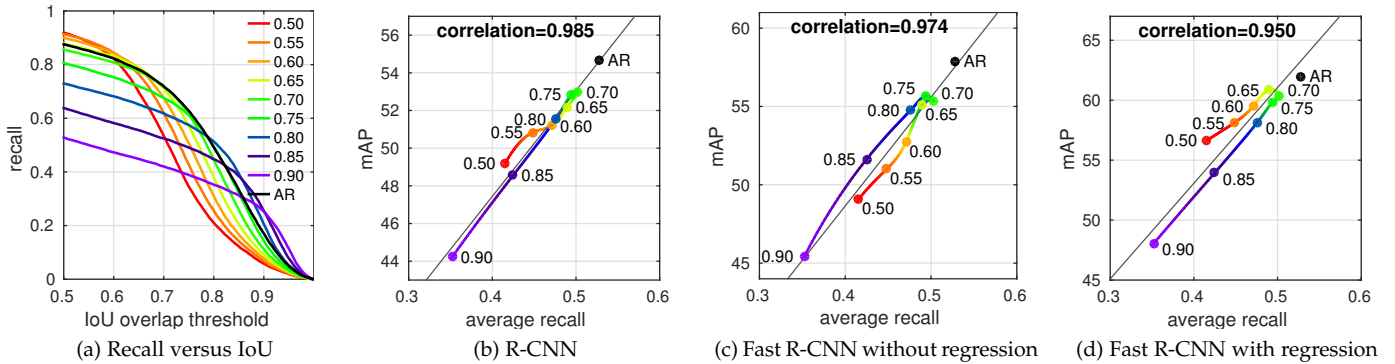


Figure 15: Finetuning EdgeBoxes to optimise AR results in top detector performance. These results further support the conclusion that AR is a good predictor for mAP and suggest that it can be used for fine-tuning proposal methods.

	aero	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
EdgeBoxesAR	69.6	78.3	66.2	58.6	42.5	82.1	78.1	83.0	42.7	74.6	66.4	81.1	82.0	74.5	68.3	35.1	66.1	68.7	75.2	62.6	67.8
+ gt oracle	75.5	79.4	70.6	63.1	55.0	82.6	84.3	83.9	46.6	75.1	68.1	82.5	83.3	75.9	76.8	41.2	67.6	70.1	77.3	65.7	71.2
+ nms oracle	77.2	87.1	76.6	67.6	48.2	84.8	85.2	87.1	52.1	83.9	72.7	86.7	87.2	84.2	77.6	44.2	75.1	73.5	83.5	65.4	75.0
+ both oracles	83.7	87.8	79.6	72.6	61.6	85.1	88.2	87.8	55.7	84.0	74.8	87.3	87.5	84.9	86.8	50.3	76.2	74.6	85.3	67.6	78.1

Table 5: Fast R-CNN (model L) detection results on PASCAL 2007 test using EdgeBoxesAR and given access to “oracles” that provide additional information to the detector. Given access to both oracles, the only way to further improve detector performance would be to avoid proposals on background or to learn a more discriminative classifier. See text for details.

adaptive. We start with a large nms threshold β_0 to encourage dense sampling around the top scoring candidates (a window is suppressed if its IoU with a higher scoring window exceeds this threshold). After selecting each proposal, β_{k+1} is decreased slightly via $\beta_{k+1} = \beta_k \cdot \eta$ to encourage greater proposal diversity. Setting $\beta_0 = 0.90$ and $\eta = 0.9996$ gave best AR at 1000 proposals on the PASCAL validation set (we kept $\alpha = 0.65$ fixed). This new adaptive EdgeBoxes variant is not optimal at any particular IoU threshold, but has best overall AR and improves Fast R-CNN mAP by 1.6 over the best previous variant (reaching 62.0 mAP).

The results in figure 15 further support our conclusion that AR is a good predictor for mAP and suggest that it can be used for fine-tuning proposal methods. We expect other methods to improve as well if optimised for AR instead of a particular IoU threshold.

5.6 Detection with oracles

We finish by exploring the limits of proposal methods when coupled with Fast R-CNN and given access to “oracles” that provide additional information to the detector. For these experiments we use the EdgeBoxesAR proposals described in §5.5 which gave the best results of all evaluated methods when coupled with the Fast R-CNN model M. Re-training the larger model L with EdgeBoxesAR proposals improves mAP to 67.8 (compared to 66.7 using SelectiveSearch proposals as in [16]).

We test two oracles. First, we augment the set of proposals with all ground truth annotations (*gt* oracle), which results in AR of 1 (but contains many false positives). Second, we perform optimal, per-class non-maximum suppression (*nms* oracle) that suppresses all false positives that overlap any true positives (without suppressing any true positives, and keeping false positives in the background). Results for the *gt* and *nms* oracles are shown in table 5.

The *gt* oracle improves mAP by about 3%. The *nms* oracle has the overall stronger effect with about 7% mAP improvement. Combining both oracles improves mAP by about 10%, indicating that their effect is largely orthogonal. All remaining mistakes that prevent perfect detection are confusions on the background or misclassifications. Therefore, the only way to further improve detector performance would be to avoid proposals on background or to learn a more discriminative classifier.

6 DISCUSSION

In this work we have revisited the majority of existing detection proposal methods, proposed new evaluation metrics, and performed an extensive and direct comparison of existing methods. Our primary goal has been to enable practitioners to make more informed decisions when considering use of detection proposals and selecting the optimal proposal method for a given scenario. Additionally, our open source benchmark will enable more complete and informative evaluations in future research on detection proposals. We conclude by summarising our key findings and suggesting avenues for future work.

Repeatability: We found that the *repeatability* of virtually all proposal methods is limited: imperceptibly small changes to an image cause a noticeable change in the set of produced proposals. Even changing a single image pixel already exhibits measurable differences in repeatability. We foresee room for improvement by using more robust superpixel (or boundary estimation) methods. However, while better repeatability for object detection would be desirable, it is not the most important property of proposals. Image independent methods such as SlidingWindow and CrackingBing have perfect repeatability but are inadequate for detection. Methods such as SelectiveSearch and EdgeBoxes seem

to strike a better balance between recall and repeatability. We suspect that high quality proposal methods that are also more repeatable would yield improved detection accuracy, however this has yet to be verified experimentally.

Localisation Accuracy: Our analysis showed that for object detection improving proposal *localisation accuracy* (improved IoU) is as important as improving recall. Indeed, we demonstrated that the popular metric of recall at IoU of 0.5 is not predictive of detection accuracy. As far as we know, our experiments are the first to demonstrate this. Proposals with high recall but at low overlap are not effective for detection.

Average Recall: To simultaneously measure both proposal recall and localisation accuracy, we report *average recall* (AR), which summarises the distribution of recall across a range of overlap thresholds. For a fixed number of proposals, AR correlates surprisingly well with detector performance (for LM-LLDA, R-CNN, and Fast R-CNN). AR proves to be an excellent predictor of detection performance both for comparing competing methods as well as tuning a specific method’s parameters. We encourage future work to report average recall (as shown in figures 7c/8c) as the primary metric for evaluating proposals for object detection. For detectors more robust to localisation errors (e.g. Fast R-CNN), the IoU range of the AR metric can be modified to better predict detector performance.

Top Methods: Amongst the evaluated methods, `SelectiveSearch`, `Rigor`, `MCG`, and `EdgeBoxes` consistently achieved top object detection performance when coupled with diverse object detectors. If fast proposals are required, `EdgeBoxes` provides a good compromise between speed and quality. Surprisingly, these top methods all achieve fairly similar detection performance even though they employ very different mechanisms for generating proposals. `SelectiveSearch` merges superpixels, `Rigor` computes multiple graph cut segmentations, `MCG` generates hierarchical segmentations, and `EdgeBoxes` scores windows based on edge content.

Generalisation: Critically, we measured no significant drop in recall when going from the 20 PASCAL categories to the 200 ImageNet categories. Moreover, while MS COCO is substantially harder and has very different statistics (more and smaller objects), relative method ordering remains mostly unchanged. These are encouraging result indicating that *current methods do indeed generalise to different unseen categories*, and as such can be considered true “objectness” methods.

Oracle Experiments: The best Fast R-CNN results reported in this paper used the large model L and `EdgeBoxesAR` proposals, achieving mAP of 67.8 on PASCAL 2007 test. Using an oracle to rectify all localisation and recall errors improved performance to 71.2 mAP, and adding an oracle for perfect non-maximum suppression further improved mAP to 78.1 (see §5.6 for details). The remaining gap of 21.9 mAP to reach perfect detection is caused by high scoring detections on the background and object misclassifications. This best case analysis for proposals that are perfectly localised shows that further improvement can only be gained by removing false positives in the proposal stage (producing fewer proposals while maintaining high AR) or training a more discriminative classifier.

Discussion: Do object proposals improve detection quality or are they just a transition technology until we have sufficient computing power? On the one hand, simply increasing the number of proposals, or using additional random proposals, may actually harm detection performance as shown in [16]. On the other hand, there is no fundamental difference between the pipeline of object proposals with a detector and a cascaded detector with two stages. Conceptually, a sliding window detector with access to the features of the proposal method may be able to perform at least as well as the cascade and as such detection proposals independent of the final classifier may eventually become unnecessary. Given enough computing power and an adequate training procedure, one might expect that a dense evaluation of CNNs could further improve performance over R-CNNs.

While in this work we have focused on object detection, object proposals have other uses. For example, they can be used to handle unknown categories at test time, or to enable weakly supervised learning [57]–[59].

Finally, we observe that current proposal methods reach high recall while using features that are not utilised by detectors such as LM-LLDA, R-CNN, and Fast R-CNN (e.g. object boundaries and superpixels). Conversely, with the exception of `Multibox` [48], none of the proposal methods use CNN features. We expect some cross-pollination will occur in this space. Indeed, there has been some very recent work in this space [60]–[62] that shows promising results.

In the future, detection proposals will surely improve in repeatability, recall, localisation accuracy, and speed. Top-down reasoning will likely play a more central role as purely bottom-up processes have difficulty generating perfect object proposals. We may also see a tighter integration between proposals and the detector, and the segmentation mask generated by many proposal methods may play a more important role during detection. One thing is clear: progress has been rapid in this young field and we expect proposal methods to evolve quickly over the coming years.

APPENDIX A ANALYSIS OF METRICS

Average recall (AR) between 0.5 and 1 can also be computed by averaging over the overlaps of each annotation gt_i with the closest matched proposal, that is integrating over the y axis of the plot instead of the x axis. Let o be the IoU overlap and $\text{recall}(o)$ the function shown for example in figure 6b. Let $\text{IoU}(gt_i)$ denote the IoU between the annotation gt_i and the closest detection proposal. We can then write:

$$\text{AR} = 2 \int_{0.5}^1 \text{recall}(o) \, do = \frac{2}{n} \sum_{i=1}^n \max(\text{IoU}(gt_i) - 0.5, 0)$$

which is the same as the *average best overlap* (ABO) [19] or the *average best spatial support* (BSS) [63] truncated at 0.5 IoU.

The ABO and BSS are typically computed by assigning the closest proposal to each annotation, i.e. a proposal can match more than one annotation. In contrast, for all our experiments we compute a bipartite matching to assign proposals to annotations (using a greedy algorithm for efficiency instead of the optimal Hungarian algorithm).

The *volume-under-surface* metric (VUS) [26] plots recall as a function of both overlap and proposal count and computes the volume under that surface. Since in practice detectors utilize a fixed number of proposals, the VUS of a proposal method is only an indirect predictor of detection accuracy.

REFERENCES

- [1] C. Papageorgiou and T. Poggio, "A trainable system for object detection," *IJCV*, 2000.
- [2] P. Viola and M. Jones, "Robust real-time face detection," in *IJCV*, 2004.
- [3] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *PAMI*, 2010.
- [4] M. Everingham, S. Eslami, L. Van Gool, C. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge – a retrospective," *IJCV*, 2014.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR*, 2009.
- [6] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft COCO: Common objects in context," *arXiv:1405.0312*, 2015.
- [7] X. Wang, M. Yang, S. Zhu, and Y. Lin, "Regionlets for generic object detection," in *ICCV*, 2013.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014.
- [9] C. Szegedy, S. Reed, D. Erhan, and D. Anguelov, "Scalable, high-quality object detection," *arXiv:1412.1441*, 2014.
- [10] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," in *ECCV*, 2014.
- [11] R. G. Cinbis, J. Verbeek, and C. Schmid, "Segmentation driven object detection with fisher vectors," in *ICCV*, 2013.
- [12] B. Alexe, T. Deselaers, and V. Ferrari, "What is an object?" in *CVPR*, 2010.
- [13] J. Carreira and C. Sminchisescu, "Constrained parametric min-cuts for automatic object segmentation," in *CVPR*, 2010.
- [14] I. Endres and D. Hoiem, "Category independent object proposals," in *ECCV*, 2010.
- [15] K. van de Sande, J. Uijlings, T. Gevers, and A. Smeulders, "Segmentation as selective search for object recognition," in *ICCV*, 2011.
- [16] R. Girshick, "Fast R-CNN," *arXiv:1504.08083*, 2015.
- [17] J. Hosang, R. Benenson, and B. Schiele, "How good are detection proposals, really?" in *BMVC*, 2014.
- [18] M.-M. Cheng, Z. Zhang, W.-Y. Lin, and P. H. S. Torr, "BING: Binarized normed gradients for objectness estimation at 300fps," in *CVPR*, 2014.
- [19] J. Carreira and C. Sminchisescu, "Cpmc: Automatic object segmentation using constrained parametric min-cuts." *PAMI*, 2012.
- [20] C. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *ECCV*, 2014.
- [21] I. Endres and D. Hoiem, "Category-independent object proposals with diverse ranking," in *PAMI*, 2014.
- [22] P. Krähenbühl and V. Koltun, "Geodesic object proposals," in *ECCV*, 2014.
- [23] P. Arbelaez, J. Pont-Tuset, J. Barron, F. Marqués, and J. Malik, "Multiscale combinatorial grouping," in *CVPR*, 2014.
- [24] B. Alexe, T. Deselaers, and V. Ferrari, "Measuring the objectness of image windows," *PAMI*, 2012.
- [25] E. Rahtu, J. Kannala, and M. Blaschko, "Learning a category independent object detection cascade," in *ICCV*, 2011.
- [26] S. Manén, M. Guillaumin, and L. Van Gool, "Prime object proposals with randomized prim's algorithm," in *ICCV*, 2013.
- [27] P. Rantalankila, J. Kannala, and E. Rahtu, "Generating object segmentation proposals using global and local search," in *CVPR*, 2014.
- [28] A. Humayun, F. Li, and J. M. Rehg, "Rigor: Recycling inference in graph cuts for generating object regions," in *CVPR*, 2014.
- [29] J. Uijlings, K. van de Sande, T. Gevers, and A. Smeulders, "Selective search for object recognition," *IJCV*, 2013.
- [30] T. Tuytelaars and K. Mikolajczyk, "Local invariant feature detectors: a survey," *Foundations and Trends in Computer Graphics and Vision*, 2008.
- [31] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. V. Gool, "A comparison of affine region detectors," *IJCV*, 2005.
- [32] T. Tuytelaars, "Dense interest points," in *CVPR*, 2010.
- [33] K. Fragkiadaki, P. Arbelaez, P. Felsen, and J. Malik, "Learning to segment moving objects in videos," in *CVPR*, 2015.
- [34] C. Gu, J. Lim, P. Arbelaez, and J. Malik, "Recognition using regions," in *CVPR*, 2009.
- [35] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *PAMI*, 2011.
- [36] P. Dollár and C. L. Zitnick, "Fast edge detection using structured forests," *PAMI*, 2015.
- [37] P. F. Felzenszwalb and D. P. Huttenlocher, "Efficient graph-based image segmentation," *IJCV*, 2004.
- [38] K.-Y. Chang, T.-L. Liu, H.-T. Chen, and S.-H. Lai, "Fusing generic objectness and visual saliency for salient object detection," in *ICCV*, 2011.
- [39] J. Lim, C. L. Zitnick, and P. Dollár, "Sketch tokens: A learned mid-level representation for contour and object detection," in *CVPR*, 2013.
- [40] M. Blaschko, J. Kannala, and E. Rahtu, "Non Maximal Suppression in Cascaded Ranking Models," in *Scandinavian Conference on Image Analysis*, 2013.
- [41] Q. Zhao, Z. Liu, and B. Yin, "Cracking BING and beyond," in *BMVC*, 2014.
- [42] P. Dollár and C. L. Zitnick, "Structured forests for fast edge detection," in *ICCV*, 2013.
- [43] J. Feng, Y. Wei, L. Tao, C. Zhang, and J. Sun, "Salient object detection by composition," in *ICCV*, 2011.
- [44] Z. Zhang, J. Warrell, and P. H. S. Torr, "Proposal generation for object detection using cascaded ranking svms," in *CVPR*, 2011.
- [45] M. Van Den Bergh, G. Roig, X. Boix, S. Manen, and L. Van Gool, "Online video seeds for temporal window objectness," in *ICCV*, 2013.
- [46] M. Van den Bergh, X. Boix, G. Roig, and L. Van Gool, "Seeds: Superpixels extracted via energy-driven sampling," *IJCV*, 2014.
- [47] J. Kim and K. Grauman, "Shape Sharing for Object Segmentation," in *ECCV*, 2012.
- [48] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, "Scalable object detection using deep neural networks," in *CVPR*, 2014.
- [49] L. Bourdev and J. Brandt, "Robust object detection via soft cascade," in *CVPR*, 2005.
- [50] H. Harzallah, F. Jurie, and C. Schmid, "Combining efficient object localization and image classification," in *ICCV*, 2009.
- [51] P. Dollár, R. Appel, and W. Kienzle, "Crosstalk cascades for frame-rate pedestrian detection," in *ECCV*, 2012.
- [52] A. Torralba, K. P. Murphy, and W. T. Freeman, "Sharing visual features for multiclass and multiview object detection," *PAMI*, 2007.
- [53] P. Zehnder, E. Koller-Meier, and L. Van Gool, "An efficient shared multi-class detection cascade," in *BMVC*, 2008.
- [54] N. Chavali, H. Agrawal, A. Mahendru, and D. Batra, "Object-proposal evaluation protocol is 'gameable'," *arXiv:1505.05836*, 2015.
- [55] D. Hoiem, Y. Chodpathumwan, and Q. Dai, "Diagnosing Error in Object Detectors," in *ECCV*, 2012.
- [56] R. Girshick and J. Malik, "Training deformable part models with decorrelated features," in *ICCV*, 2013.
- [57] S. Vicente, C. Rother, and V. Kolmogorov, "Object cosegmentation," in *CVPR*, 2011.
- [58] M. Guillaumin, D. Kuttel, and V. Ferrari, "Imagenet auto-annotation with segmentation propagation," *IJCV*, 2014.
- [59] K. Tang, A. Joulin, L.-J. Li, and L. Fei-Fei, "Co-localization in real-world images," in *CVPR*, 2014.
- [60] W. Kuo, B. Hariharan, and J. Malik, "Deepbox: Learning objectness with convolutional networks," *arXiv:1505.02146*, 2015.
- [61] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *arXiv:1506.01497*, 2015.
- [62] P. O. Pinheiro, R. Collobert, and P. Dollár, "Learning to segment object candidates," *arXiv:1506.06204*, 2015.
- [63] T. Malisiewicz and A. A. Efros, "Improving spatial support for objects via multiple segmentations," in *BMVC*, 2007.