

---

# SLOT-BASED IMAGE AUGMENTATION SYSTEM FOR OBJECT DETECTION

---

A PREPRINT

**Yingwei Zhou\***

School of Electronics and Computer Science  
University of Southampton  
Southampton, SO17 1BJ  
yz39g15@soton.ac.uk

## ABSTRACT

Object Detection has been the significant topic in computer vision. As the continuous development of Deep Learning, many advanced academic and industrial outcomes are established on localising and classifying the target objects, such as instance segmentation, video tracking and robotic vision. As the core concept of Deep Learning, Deep Neural Networks (DNNs) and associated training are highly integrated with task-driven modelling, having great effects on accurate detection. The main focus of improving detection performance is proposing DNNs with extra layers and novel topological connections to extract the desired features from input data. However, training these models can be a computational expensive and laborious progress as the complicated model architecture and enormous parameters. Besides, dataset is another reason causing this issue and low detection accuracy, because of insufficient data samples or difficult instances. To address these training difficulties, this thesis presents two different approaches to improve the detection performance in the relatively light-weight way. As the intrinsic feature of data-driven in deep learning, the first approach is "slot-based image augmentation" to enrich the dataset with extra foreground and background combinations. Instead of the commonly used image flipping method, the proposed system achieved similar mAP improvement with less extra images which decreases training time. This proposed augmentation system has extra flexibility adapting to various scenarios and the performance-driven analysis provides an alternative aspect of conducting image augmentation.

**Keywords** Image Augmentation · Object Detection · Deep Learning · Computer Vision

## 1 Introduction

Object detection is a fundamental research area to answer where are the objects and what they are. It consists of two key sub-problems of localisation and classification. Despite its simple conception, many advanced computer vision tasks are based on precisely acquired object location and categories by object detection methods, such as instance segmentation, video tracking and autonomous driving [1, 2, 3]. Therefore detection quality has a great impact on these advanced vision tasks. The main contributions of this thesis are introducing two novel aspects to improve the detection performance which are light-weight, flexible to extend and less computational cost. To be specific, an image augmentation system called "slot-based image augmentation" is proposed and experiments have shown its improvements for small object detection and huge potential capacity in augmenting images. Moreover, it is implemented in an image pre-processing fashion, no need for training and easy to extend for various scenarios. This approach highlights the significance of data-source, addressing the very intrinsic of data-driven machine learning approaches. Besides, this work has found failure of commonly used image augmentation methods that they might damage the mAP of detection DNNs.

In practice, the limitation of available data is a common case and data augmentation methods such as flipping, rotating and tuning brightness, are applied to get through the bottleneck of lacking training data. These methods are conducted

---

\*GitHub: <https://mercurise.github.io/>

as manual image transformation without task-oriented features. Although these methods contribute partial increase for image classification accuracy, they are not that beneficial for object detection. In some cases, augmentation makes the model performance worse and related results are presented in this Chapter. The main reason of the failures is the differences in target learn-able features between image classification and object detection. In addition, object detection requires extra location and contextual features rather than the object feature alone as in classification task. Moreover the manual image transformation based augmentation are not guaranteed to augment sufficient extra learn-able features especially for small scale objects.

Therefore, the above discussion draws the motivation that an image augmentation approach is desired specially for object detection task to provide extra classification and localisation features. Regarding the methodological gap between object detection and regular augmentation methods, the slot-based image augmentation is proposed to generate images with more learn-able object detection related features by adding extra combinations of foreground objects and background images.

## 2 Related Works

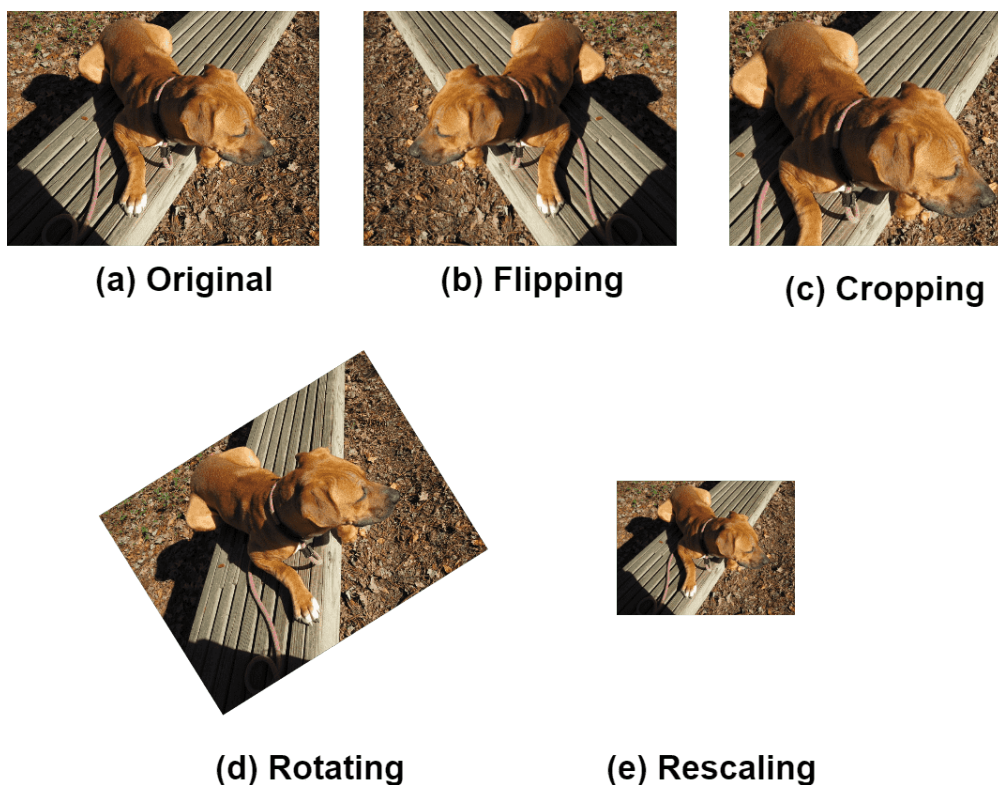


Figure 1: The Examples of Commonly Used Image Augmentation Methods.

Commonly used image augmentation approaches include manual image transformations such as flipping, rotating, cropping and re-scaling as shown in Figure 1. Because images are represented as matrix in computer vision, these transformations are conducted by linear transformation on matrix and many toolboxes support these transformations such as OpenCV [4], Scikit-image [5] and recently released “Population Based Augmentation” by [6]. Besides of these transformations on image size or positions, another alternative approach is the colour-wise augmentation including brightness, contrast and saturation, which are supported by many deep learning frameworks: TorchVision [7], mxNet [8], etc. In summary, most of the augmentation works are based on simple image transformation and there are not many other augmentation approaches featured with performance analysis or adaptive augmentation.

Besides these augmentation methods, there are some works on inserting objects into background image to generate new images. [9] proposed a machine learning based method to search suitable positions to insert the objects. These positions are ranked by the performance evaluation of the applied object detection models.

Among these works, [10] highlights the importance of scales and [11] conducted similar works on small object detection.

Instead of these augmentation related works, the main stream focuses on modifying DNN architecture to improve detection performance by mixing CNN feature maps of different layers. For the small object detection, the contextual information is highlighted thus various DNN architectures are proposed to utilise these features as possible such as original feature pyramid pooling by [12], the relation based DNN architecture by [13] and contextual refinement methods by [14].

These architectural modifications methods achieved desired detection improvements however, added extra computation costs and less flexible for transferring the trained models to different scenarios. Inspired by these proposed novel methods and impressive implementation works, the proposed slot-based image augmentation method is lightweight, flexible for transfer learning and highlights a general way of augmenting images for object detection.

### 3 slot-Based Image Augmentation System

#### 3.1 System Introduction

The slot-based augmentation system is built on the fundamental element called "slot" which is a generalised conception from the isolated foreground objects. In other words, a slot is the replaceable position which is initially the isolated foreground object in an image and the slot-based image augmentation produces extra images by substituting foregrounds enriching various combinations of foreground and backgrounds. Furthermore detecting objects is the process of recognising foreground objects from its background utilisation the features learnt from training images. Therefore the slot substitution enriches the combinations especially for the unbalanced dataset in which there are not sufficient objects of a certain category. In addition, the slot based augmentation is easy to customise and highly flexible for different scenarios such as substituting objects of the same category, which creates extra learn-able "foreground and background" correlation features for the target category objects. Figure 2 example of augmentation. It is also a potential way to change the bounding box style slot substitution to an instance wise substitution as the instances for segmentation tasks are represented as polygons.

## 4 System Implementation

Regarding the conception of slots, the two highlighted features are foregrounds and isolation. On one hand, foregrounds assign slot location as the foregrounds which tightly include objects inside the bounding box and causes limited corruption for other objects or background while substituting them. On the other hand, the isolation selects individual objects that have no insertion with any other objects and no damage to other objects when replacing them with other ones. Hence a slot is a rectangle area of the background image with the same position as an isolated original object and causes minor effects while substitute it.

Despite the simple definition, it takes several steps to conduct the augmentation and the system is designed in an end-to-end fashion making it flexible to modify and extend with more functionalities, which are presented in the following section.

#### 4.1 System Design and Pipeline

As introduced above, slot is the central component of the augmentation system which is required to be isolated and possessing the foreground object locations. Dataset for object detection includes images and annotation files that specify object location as coordinates and object information such as category and scales. Thus foreground object locations are capable to be directly extracted from the annotation files while isolation requires extra steps to extract them. The progress of selecting isolated objects is simplified as finding individual rectangle bounding boxes that have no insertion with any others in the same image. As bounding boxes are represented with coordinates, a heuristic method is applied to solve this problem by comparing a set of coordinates as in Figure and the rules are shown below:

For two bounding boxes represented by the coordinates of the top-left and bottom-right vertexes as  $bbox1: [x_{11}, y_{11}, x_{12}, y_{12}]$  and  $bbox2: [x_{21}, y_{21}, x_{22}, y_{22}]$ , where  $(x_{11}, y_{11})$  and  $x_{21}, y_{21}$  are top-left vertexes while the bottom-right ones are  $x_{12}, y_{12}$  and  $x_{22}, y_{22}$ . So  $bbox1$  and  $bbox2$  are overlapped if and only if satisfying all the following conditions and the isolated slots are decided by the complementary cases:

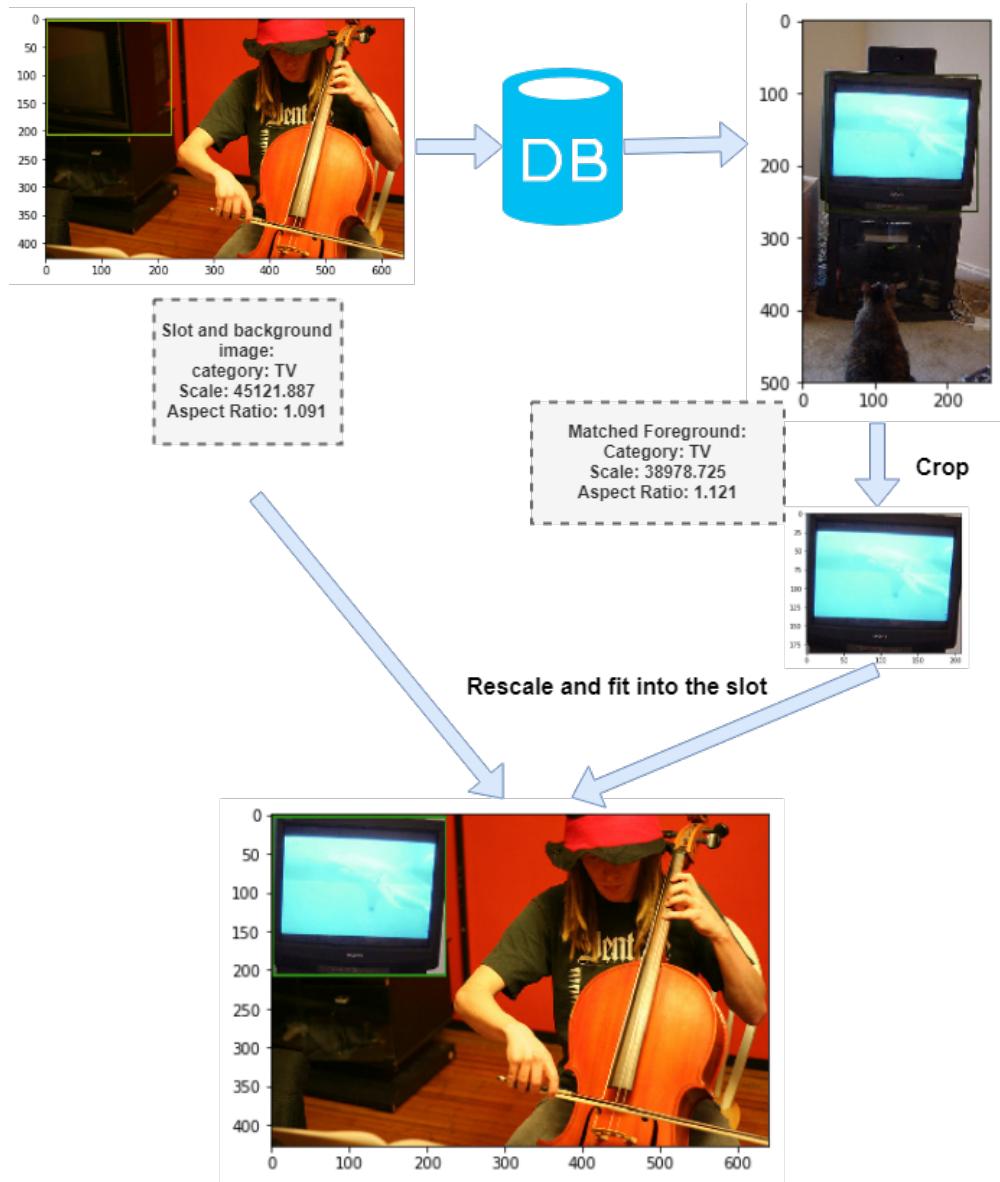


Figure 2: Generating one image by fitting one TV slot with a valid foreground object which has the similar aspect ratio and scales.

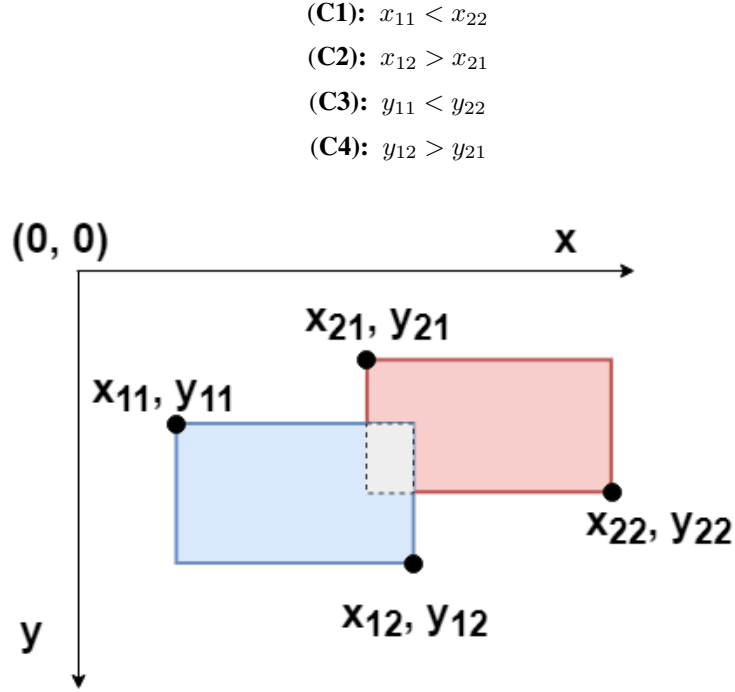


Figure 3: Slots are the isolated foregrounds represented as rectangle bounding boxes. Finding isolated bounding boxes can be achieved by analysing coordinates according to a series rules.

The system pipeline contains two main phases: initialisation and image augmentation as shown in Figure 4. To be specific, the valid slots selected by matching the complementary conditions proposed above. These picked slots are recorded into a database with detailed information include slot location, original foreground category and other scale related values (width, height, area size and aspect ratio). Generating images is conducted by substituting slots satisfying specific schemes determined by certain scenarios and scenarios are determined based the model performance evaluated by mAP and Derek P-R curve. For example, if the model behaves low mAP at a certain category, the augmentation is set to produce images containing objects of the target category, Figure 2 gives a detailed workflow of augmenting extra images containing objects of the "television" category. Many other scenario augmentation is presented in later sections of this chapter.

In the reproduction and implementation aspect, the system consists several central steps such as system initialisation, performance analysis of detection model and image augmentation based on the decided scenario. In experiments, MS-COCO 2014 dataset is selected as the main dataset as it provides helpful API codes and contains a large amount of objects especially small objects [15]. MS-COCO 2014 is split as full training set, a subset of validation called "val minus minival" and the reduced validation dataset called "minival". The training set and randomly selected 5k "val minus minival" are used for training while "minival" is used for validation. Furthermore, MS-COCO is a commonly used dataset in many related works as the reference of model capacity. For detection model selection, the PyTorch ([7]) version faster RCNN is selected as the benchmark detection model implemented by [16] which uses ImageNet pre-trained ResNet-50 as its backend proposed by [17, 18]. Faster RCNN is a widely used model to compare with related works or transfer pre-trained weights and sufficiently supported by the deep learning community.

As in the pipeline Figure 4, system initialisation follows loading dataset into the system, select isolated slots and summarise them into a database. In practice, dataset loading is implemented based on customised MS-COCO API to read annotation (JavaScript), load training images and refine the data structure. The associated Python libraries include Numpy, Scikit-Image, Matplotlib and many others [19, 5, 20]. The isolated slots selection is conducted by implementing the condition matching mentioned above, mainly by Numpy library. Database is implemented using data structure of Python list and dictionary encapsulated by PANDAS library ([21]), which is feasible for matching SQL (Structured Query Language) style enquiry in a fast and efficient way especially handling large amount of data records. In addition, the slot aspect ratio, area size and other information are calculated using Numpy and Jupyter Notebook is the principle tool for efficient coding.

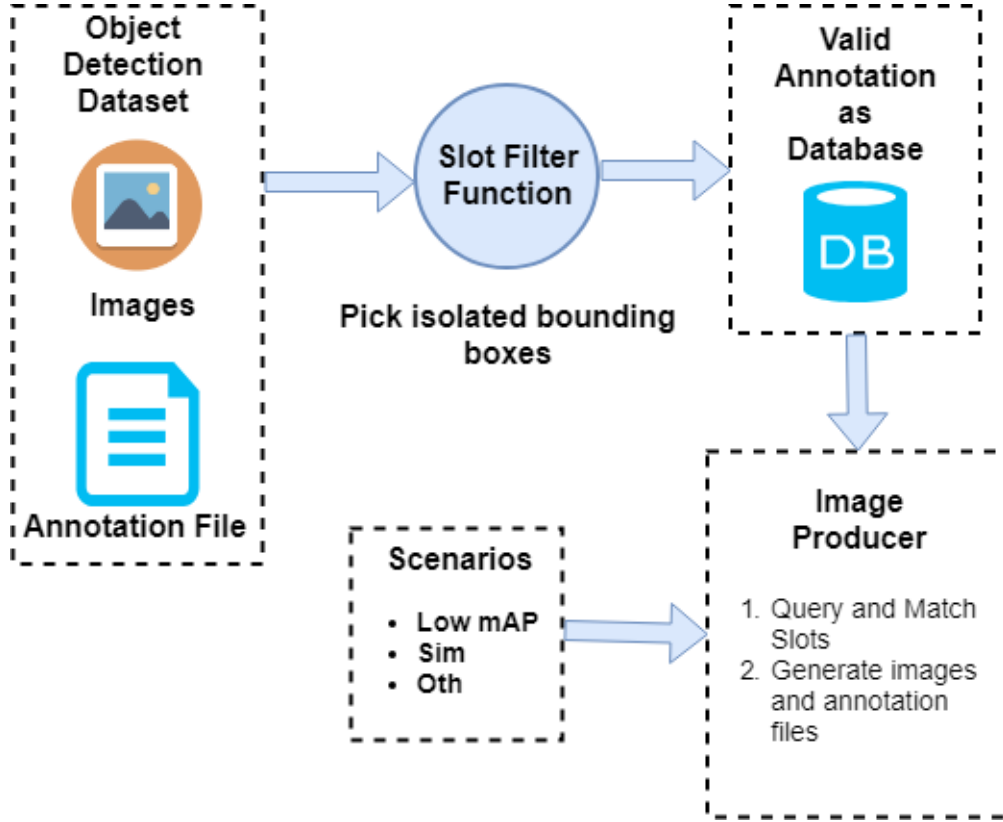


Figure 4: Components of the Slot-based Image Augmentation System.

To analyse model performance with mAP, MS-COCO API provide evaluation tools in Python, Matlab and other languages. In practice, Python version official evaluation tools are integrated in the implementation of faster RCNN benchmark model. Since Derek P-R curve plotting is not supported by Python API, the Matlab version is used for the curve drawing and extra Linux shell scripts are implemented for the intermediate data structural transformation between Python and Matlab API.

Finally for the slot substitutions, the slot based sub-image cropping and pasting is conducted by python libraries of OpenCV and Scikit-Image [4, 5]. Besides the crop and paste step, there are normally multiple available candidates for a certain slot and the rules of candidate selection is another important problem to concern.

## 5 Filtering Slot Candidates

As mentioned in previous sections, slot-based image augmentation is conducted by substituting foreground objects and re-scale the candidate foreground to fit the slot shape. It is a common case that many candidate foregrounds are suitable for the target slot. Hence a candidate selection rule is required to address the “many-to-one” issue. In this section, three different filters are introduced to ensure the augmentation quality with analysis and experiments presented.

In practice, an “attribute matching” strategy is proposed to select valid candidates for the substitution with less damage to the potential learn-able features of the background image and candidate. Regarding the instance information in a large dataset, instances have a wide range of scale related attributes such as area size, length, width and height. Among those attributes, the instance area size, aspect ratio and category are applied as the filters to select candidates for slot substitution, in which the area size and aspect ratio are calculated in the format of rectangle bounding boxes.

As the example in Figure 5, candidates are firstly filtered by their aspect ratio to avoid invalid re-scaling. Then candidates with similar scales/resolution are selected by the scale filter. Final step is filter rest of these candidates by their categories. The category filter is a scenario driven component. For example, candidates of the same category as the original slot foreground are selected when the augmentation system is aimed at augmenting extra images for a certain

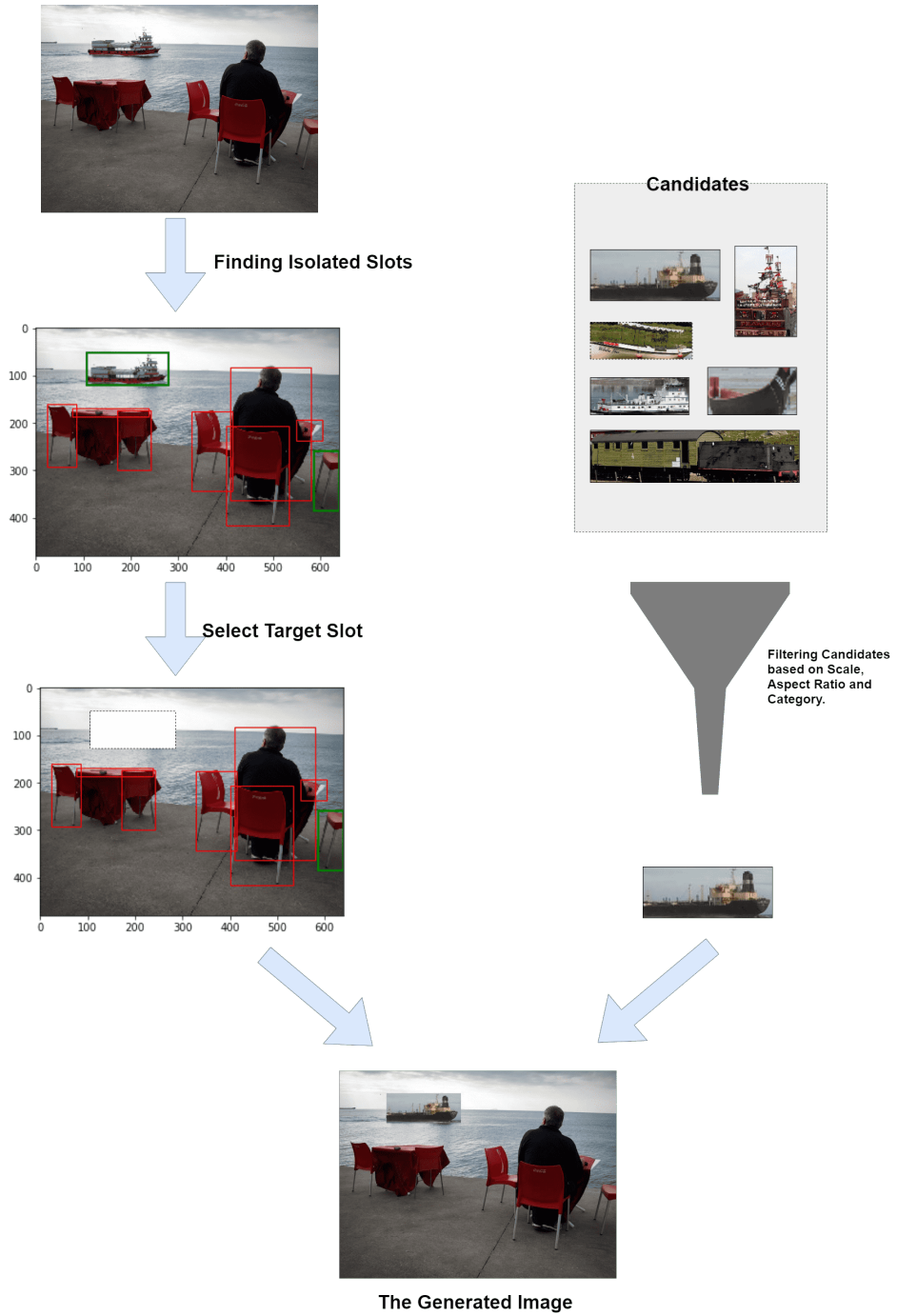


Figure 5: System Work-flow



category. Another scenario is augmenting images to improve detection performance of one certain super-category instances such as substitute cats with dogs under animal super-category. Overall, filters are decided to preserve as much object information as possible without corrupting background images. These three filters are simple, effective and easy for modification. Candidates selected by the three filters are randomly chosen as the final candidate to fit the target slot.

Furthermore, several experiments are conducted to test the validity of selecting these filters and a loss function based candidate matching policy is introduced to compare with the random selection method for the final candidate decision. To transfer the scenario in a more realistic level, a mini-COCO dataset is established for filter validation tests, presented in the following section.

### 5.1 Filter Validation on “mini-COCO” dataset

Table 1: Top and Bottom Three Categories of Instance Amount in MS-COCO

Categories	Image Amount	Instance Amount
<b>Person</b>	<b>64,115</b>	<b>262,465</b>
Chair	12,774	38,491
Car	12,251	43,867
....		
Parking Meter	705	1,285
Toaster	217	225
<b>Hair Drier</b>	<b>189</b>	<b>198</b>

As previously mentioned, MS-COCO dataset consists of images and instances, providing rich resource for slots and substitution candidates. However, lacking of data samples is a common issue due to the limited amounts of slots and candidates. To ensure the selected slot-based filters are functional in general conditions, it is necessary to conduct validation experiments on simulated scenarios. Hence, a mini-COCO dataset is established to run these tests with reduced image and instance amounts, limited slots and preserving the consistence with original MS-COCO. The mini-COCO dataset is required with the followings features:

1. Reduced instance and image amounts to simulate the realistic scenario.
2. All MS-COCO instance categories/super-category are included.
3. Reasonable slot amount.
4. Relatively balanced in an instance level. Table 5.1 has listed some of the category amounts in MS-COCO, in which some categories have much more instances than the other.

To maintain those desired features, mini-COCO dataset is established in a heuristic way to iteratively fetch images from MS-COCO as presented in Algorithm 5.1. Based on MS-COCO annotations, all 80 categories are sorted in an ascending order of instance amount. After some other data structure wise setup, the desired mini-coco dataset is selected from a set of candidates delivered by an accumulative operation of stacking images containing a certain category instances. In other words, a single category with the smallest instance amount is selected by previous sorting operation and those images containing instances of the selected category, are added into mini-coco as a candidate. As the iterative progress, these mini-coco candidates are evaluated with a series of metrics including standard deviation of instance amounts for dataset balancing concern, slot amount, summary of instance and image amounts and whether all categories are included. The final mini-coco is selected based on the comparison of those feature-related metrics and in our case, mini-coco dataset is selected at the 20th category is selected containing 73,653 instances in 10,436 images with all instance categories included. In addition, mini-coco contains approximately 9% images (8.8%) and instances (8.6%) comparing with the original MS-COCO which contains 118,287 images and 860,001 instances. Figure 6 has presented the related metric values as stacking categories into the mini-coco candidates.

### 5.2 Filter Details and Validation

As discussed above, the chosen filters are aspect ratio, scale and category. In this section experiments are conducted on testing turning on and off these filters using mini-coco dataset and baseline faster RCNN model.



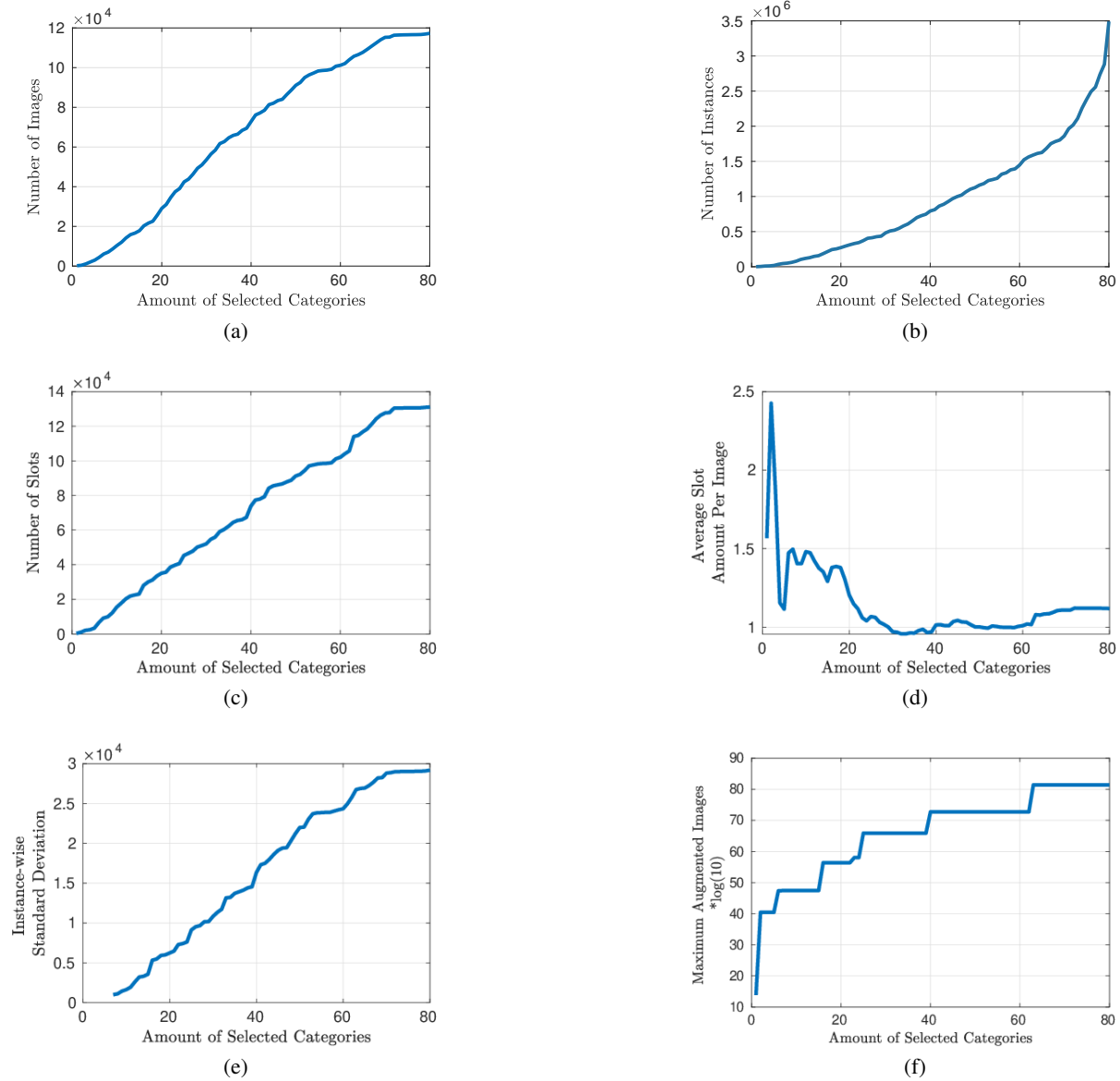


Figure 6: **Row one:** image and instance amount as selecting more categories. **Row two:** overall slot amount and average slot number per image of mini-coco as selecting more categories. **Row three:** instance standard deviation and augmentation capacity.

---

**Algorithm 1** A heuristic algorithm to construct mini-COCO dataset

---

```

1: procedure CONSTRUCT MINI-COCO
2:   Initialisation:
3:   num_cat  $\triangleright$  category amount in MS-COCO
4:   sorted_array  $\leftarrow$  MS-COCO categories sorted by instance amount in an ascending order
5:   mini_coco  $\leftarrow$  []
6:   records  $\leftarrow$  []  $\triangleright$  calculate and store values corresponding to dataset requirements including standard deviation
   (std), slot amount (slot_amount) and check whether all categories are included (category_included)
7:   Progress
8:   while i from 0 to num_cat - 1 do  $\triangleright$  iterate items in sorted_array
9:     mini_coco[i]  $\leftarrow$  images and associated annotations containing category sorted_array[i]
10:    CLEANUP  $\triangleright$  remove selected images and annotations from MS-COCO
11:    records[i]  $\leftarrow$  std, slot_amount and category_included
12:   Output records and mini_coco
13:   RETURN mini_coco[a]  $\triangleright$  mini_coco[a] is manually chosen based on records and mini_coco.

```

---

Table 2: Filter ON/OFF Validations

Category	Scale	Aspect Ratio	mAP
ON	ON	ON	<b>0.149</b>
OFF	ON	ON	0.146
ON	ON	OFF	0.147
ON	OFF	ON	0.148

In practice, filters are defined as a selector to pick valid candidates from a certain range w.r.t. the scheme. For example, the scale filter compares candidates’ scales and preserve those within  $\pm 20\%$  of the original scales. The aspect ratio filter works in the same way and the category filter. The proportion threshold  $\pm 20\%$  is a trade off between candidate quality and candidate amount. The threshold is preferred to be larger When the dataset includes large number of instances and smaller with relatively small dataset. Category filter is applied according to specific scenario. For the task of generating more images with same category, candidates are preserved only if they belong to the same category as the original slot instances.

Figure 5.2 shows the experiments on baseline model testing turning on or off those filters on mini-coco dataset. The experimental results demonstrates the necessary of turning on these filters and Figure 7 gives an invalid example image without turning on the scale filter and Figure 8 presents the case of inappropriate aspect ratio.

## 6 Performance-based Image Augmentation

With the aforementioned filters, slot-based augmentation system can be applied in many different scenarios and this section presents an example of improving detection performance by generating images with the substitution of same category instances.

The performance-based image augmentation is based on the training mAP and P-R details of the baseline model. In this case, the original faster RCNN model is selected as baseline model and trained on MS-COCO dataset with 490k batches without any augmentation (same hyperparameter configuration as [22]). The baseline model produces 30.5% overall mAP (C75) that reaches the original benchmark. However, mAP alone is biased to describe model performance and weakness. Hence, detailed P-R curves in Figure 9 are introduced to evaluate the scale-wise performance. In addition to the overall performance, category-wise P-R curves are analysed as well. The P-R curve and detailed mAP performance indicates the imbalanced performance that some categories have fairly low mAP due to limited image and instance amount. To address the issue of low mAP of certain categories, augmenting extra images to provide extra learn-able features is expected a performance improvement, hence the motivation of using slot-based augmentation system.

Based on the analysis of the mAP and instance amount, three categories are selected to demonstrate the process of augmenting images to improve the performance: **cars**, **ships** and **traffic lights**. These three categories initially have less instance amounts, lower detection mAP and reasonable number of slot. To augment more images containing these

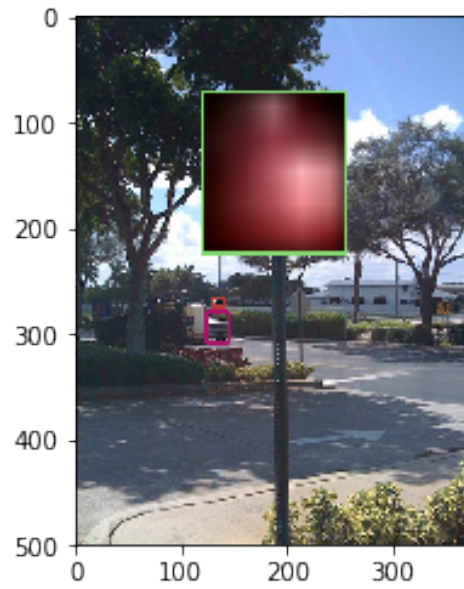


Figure 7: A negatively generated image, in which the street sign is too small to fit into the slot.

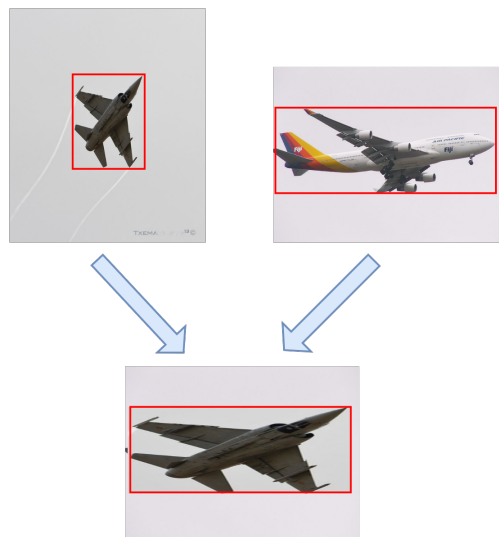


Figure 8: Object feature loss is caused by substituting the target slot with an instance of different aspect ratio. In this case, the target slot ratio is width larger than height while the substitution is opposite.

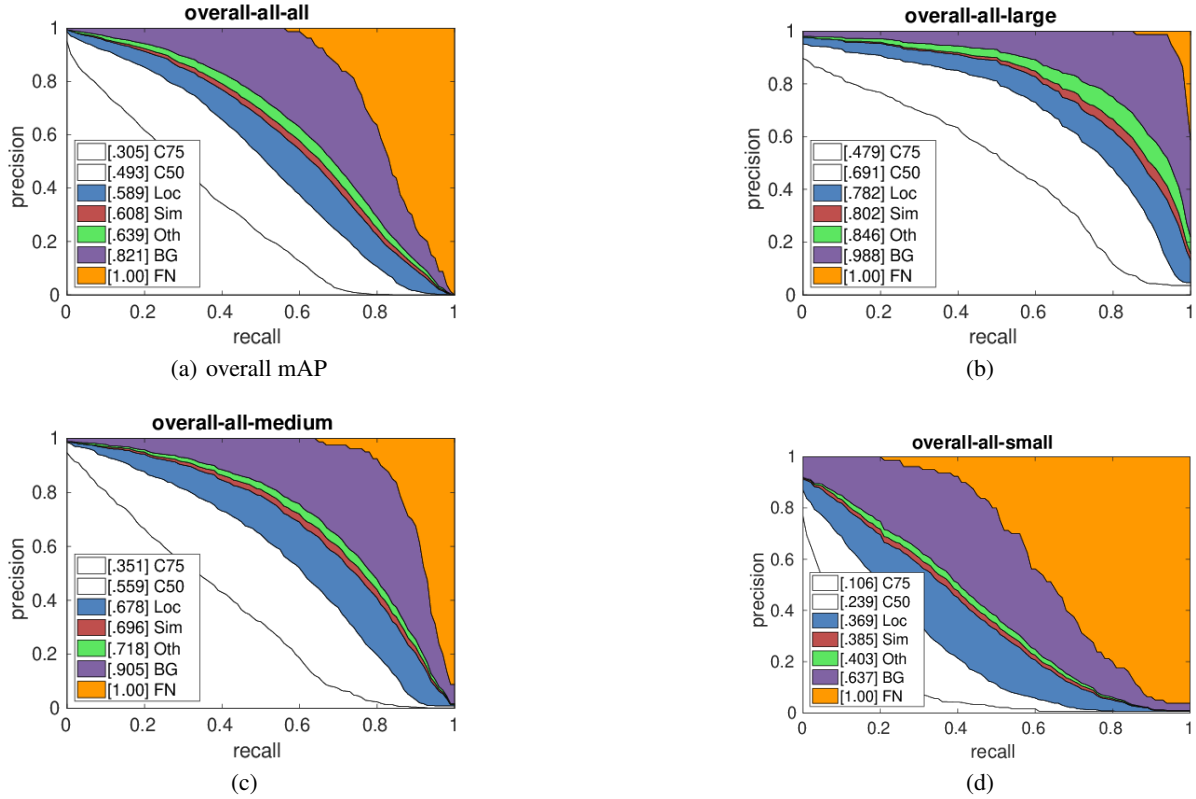


Figure 9: Derek Style P-R Curves of Baseline Faster RCNN Model without Image Augmentation.

Table 3: Number of Augmented Images with Different Approaches

Augmentation Methods(490k)	Number of Original Images	Number of Augmented Image	Proportion of Original Images
No-Flipping	118,287	0	0%
Flipping (All Categories)	118,287	118,287	50%
Ours Cars	12,251	3,262	78.97%
Ours Boats	3,025	940	76.29%
Ours Traffic Lights	4,139	2,224	64.05%
Flipping and Ours Cars	12,251	15,513	44.13%
Flipping and Ours Boats	3,025	3,956	43.31%
Flipping and Ours Traffic Lights	4,139	6,363	39.41%

category instances, the slots are substituted with candidates of the same category and the  $\pm 20\%$  threshold is applied on the scale and aspect ratio filters. Based on the selected categories, the widely used flipping and slot-based approaches are compared in the aspects of various mAP metrics, the amount of generated images and the effects of combining them together. Table 6 displays the augmentation amount with different augmentation methods, in which flipping images is the commonly used approach in many detection systems. In addition, these augmented images also change the probability of training on an original image in other words, the chances of learning original features or artificial features. Details of the augmentation performance are presented in the following sections. **Note** the slot-based method is conducted one epoch so that one slot is only substituted with one candidate.

### 6.1 Augmenting Cars Images

MS-COCO contains 43,867 car instances distributed in 12,251 images. Following the pipeline discussed in Section 4.1, isolated slots are selected and replaced with candidates satisfying the requirements of the filters. In other words, the car

slots are substituted with those candidates with close aspect ratio, scales and the same category. In the experiments, 3,262 images are generated by replacing slots by one iteration and the sample images are shown in Figure 12(a).

Regarding augmentation effects on detection performance, Table 6.1 presents baseline model mAP with different augmentation methods of all scales. Table 6.1, Table 6.1 and Table 6.1 describe the same schemes on objects with scales of large, medium and small relatively. To summarise these results, slot-based augmentation system generates around 3000 images (over 12,251 images) and train on these extra images improved the overall mAP 2% comparing with commonly used “flipping” and non-augmentation. Regarding the training mAP of different scales, it can be observed that slot-based augmentation has greater effects for medium and small objects, boosting 1.6% mAP for medium instances and 0.9% for small scale (comparing with not applying any augmentation methods, under the C75 metric). Furthermore, the mixing of slot-based method and flipping, increased 1.4% mAP while 0.5% for flipping only. Despite the mAP improvements of flipping, the overall mAP from Table 6.1 shows a 0.7% mAP decreasing which might be the reason of medium scale instances in Table 6.1.



Figure 10: **Row one:** augmented car images. **Row two:** the original images (some images are cropped for presentation purpose).

Table 4: Baseline Model mAP on Cars

	C75	C50	Loc (C10)
No Flipping	0.268	0.513	0.688
Flipping	0.261	0.510	0.691
<b>Ours</b>	<b>0.281</b>	<b>0.509</b>	<b>0.679</b>
<b>Flipping + Ours</b>	<b>0.276</b>	<b>0.515</b>	<b>0.677</b>

Table 5: Baseline Model mAP on Cars (Large Scale)

	C75	C50	Loc (C10)
No Flipping	0.538	0.776	0.874
Flipping	0.599	0.768	0.872
<b>Ours</b>	<b>0.580</b>	<b>0.778</b>	<b>0.883</b>
<b>Flipping + Ours</b>	<b>0.538</b>	<b>0.802</b>	<b>0.883</b>

Table 6: Baseline Model mAP on Cars (Medium Scale)

	C75	C50	Loc (C10)
No Flipping	0.497	0.727	0.860
Flipping	0.486	0.727	0.842
<b>Ours</b>	<b>0.513</b>	<b>0.736</b>	<b>0.848</b>
<b>Flipping + Ours</b>	<b>0.501</b>	<b>0.729</b>	<b>0.841</b>

Table 7: Baseline Model mAP on Cars (Small Scale)

	C75	C50	Loc (C10)
No Flipping	0.118	0.384	0.620
Flipping	0.123	0.372	0.622
<b>Ours</b>	<b>0.127</b>	<b>0.374</b>	<b>0.611</b>
<b>Flipping + Ours</b>	<b>0.132</b>	<b>0.376</b>	<b>0.599</b>

## 6.2 Augmenting Boats Images

Similar to augment images with car instances, after one iteration, 940 boat images are generated from 3,025 images with 10,759 instances. Figure 11 shows three sample images before and after slot substitution. Considering the training mAP, the flipping and slot-based augmentation methods actually decreased the mAP by 0.9 % and 3.1% respectively. While combining these two augmentation methods yields an increasing 2.3% mAP. For different scale objects, the two augmentation methods decreased detection mAP at some extent. But slot-based method contributes a 2.8% improvement for medium scale objects and flipping method rises a 1.3% improvement on small scale objects. In general, for the boat-related image augmentation, individual augmentation failed to improve training mAP with a decreasing effect instead. However, applying these two methods together improved detection mAP at all scales.



Figure 11: **Row one:** augmented boat images **Row two:** the original images

## 6.3 Augmenting Traffic Light Images

For the traffic light category, 2,224 images are augmented based on 12,884 instances from 4,139 original images within one epoch. Some sample images are shown in Figure 12(a). In the aspect of detection mAP, both flipping and slot-based

Table 8: Baseline Model mAP on **Boats**

	C75	C50	Loc (C10)
No Flipping	0.127	0.37	0.594
Flipping	0.118	0.366	0.600
<b>Ours</b>	<b>0.096</b>	<b>0.351</b>	<b>0.592</b>
<b>Flipping + Ours</b>	<b>0.140</b>	<b>0.385</b>	<b>0.603</b>

Table 9: Baseline Model mAP on **Boats (Large Scale)**

	C75	C50	Loc (C10)
No Flipping	0.342	0.639	0.778
Flipping	0.323	0.591	0.769
<b>Ours</b>	<b>0.239</b>	<b>0.570</b>	<b>0.784</b>
<b>Flipping + Ours</b>	<b>0.428</b>	<b>0.630</b>	<b>0.795</b>

Table 10: Baseline Model mAP on **Boats (Medium Scale)**

	C75	C50	Loc (C10)
No Flipping	0.125	0.474	0.727
Flipping	0.103	0.460	0.728
<b>Ours</b>	<b>0.153</b>	<b>0.504</b>	<b>0.757</b>
<b>Flipping + Ours</b>	<b>0.134</b>	<b>0.428</b>	<b>0.711</b>

Table 11: Baseline Model mAP on **Boats (Small Scale)**

	C75	C50	Loc (C10)
No Flipping	0.043	0.259	0.564
Flipping	0.056	0.272	0.585
<b>Ours</b>	<b>0.029</b>	<b>0.230</b>	<b>0.563</b>
<b>Flipping + Ours</b>	<b>0.063</b>	<b>0.275</b>	<b>0.565</b>

augmentation have improved the mAP. The improvement for flipping is 0.3% and 2.0% for slot-based method. The combination of these two method is 2.0% as well. Observing the mAP changes on different scales, flipping method improved 3.2% and 4.9% for large and medium objects while a decreased 0.4% for small objects. The slot-based method produces 2.9% and 3% for small objects. However, it decreased 0.7% for medium objects. The combination approach improved all the three scales with 3.5%, 3.3% and 1.5% for large, medium and small objects respectively.

Table 12: Baseline Model mAP on **Traffic Lights**

	C75	C50	Loc (C10)
No Flipping	0.096	0.363	0.539
Flipping	0.099	0.366	0.553
<b>Ours</b>	<b>0.116</b>	<b>0.377</b>	<b>0.534</b>
<b>Flipping + Ours</b>	<b>0.116</b>	<b>0.379</b>	<b>0.547</b>

Table 13: Baseline Model mAP on **Traffic Lights (Large Scale)**

	C75	C50	Loc (C10)
No Flipping	0.431	0.658	0.865
Flipping	0.463	0.687	0.853
<b>Ours</b>	<b>0.460</b>	<b>0.656</b>	<b>0.833</b>
<b>Flipping + Ours</b>	<b>0.566</b>	<b>0.727</b>	<b>0.896</b>

## 7 Discussion and Analysis

As the main metrics of evaluating object detection models, mAP provides an overview of model performance. However, for different category instances, the mAP per category varies a lot. A normal case is that, categories with large



Figure 12: **Row one:** augmented traffic-light images **Row two:** the original imagesTable 14: Baseline Model mAP on **Traffic Lights (Medium Scale)**

	C75	C50	Loc (C10)
No Flipping	0.265	0.683	0.811
Flipping	0.316	0.682	0.823
<b>Ours</b>	<b>0.258</b>	<b>0.660</b>	<b>0.837</b>
<b>Flipping + Ours</b>	<b>0.298</b>	<b>0.686</b>	<b>0.825</b>

Table 15: Baseline Model mAP on **Traffic Lights (Small Scale)**

	C75	C50	Loc (C10)
No Flipping	0.050	0.291	0.478
Flipping	0.046	0.299	0.497
<b>Ours</b>	<b>0.080</b>	<b>0.316</b>	<b>0.466</b>
<b>Flipping + Ours</b>	<b>0.065</b>	<b>0.311</b>	<b>0.485</b>

amount of instances comes with higher mAP such as the “person” category with 44.4% mAP (C75) containing 262,645 instances in 64,115 images. While those with small instance amounts have much lower mAP such as “hair drier” 1% mAP (C75) with only 198 instances in 189 images. Previous section has presented the experiments of applying slot-based augmentation method into three different category instances. These experimental results have shown the benefits of applying augmentation methods. In general, applying image augmentation is expected a higher overall detection mAP. However, the augmentation affects differently over categories and sometimes it even decreased the detection mAP, such as applying flipping method in boat objects. Compared with flipping images, slot-based augmentation produced closing mAP improvement. For example, the augmentation on car objects behaved even better than flipping with an extra 2% mAP. Besides, the slot-based augmentation added less images comparing with flipping which doubled the image amount. For a large dataset such as MS-COCO, flipping actually increased a large number of images causing extra training time and computation resources. While slot-based augmentation method is relatively more flexible to control the amount of augmentation with less increased computation time. As the performance for some categories is not improved, combining these methods together produced a consistent improvement over all these three categories. For the large boat instance, both flipping and slot-based augmentation failed to improve mAP while the combination provides a 8.6% increase.

In conclusion, image augmentation is frequently applied in object detection based tasks to provide extra learn-able features and improve the detection mAP. One of the commonly used method is flipping all the images of the dataset. In this Chapter, a slot-based image augmentation method is proposed, in which images are augmented by replacing isolated foregrounds to provide extra combinations of foreground and backgrounds. Additionally, the system components are tested by a series of filter experiments and an augmentation scenario is presented showing detailed procedure of applying this method. Besides the expected mAP improvement with augmentation methods, the experimental results have highlighted several interesting facts:

- Instance amounts affects category-wise mAPs. Some categories with lower mAP is resulted from lacking of data, which highlights the importance of applying image augmentation methods on those specific images without rising extra training time. In addition, detailed mAP and the Derek P-R curve is descriptive on analysing detailed performance.
- Image augmentation is not always bringing benefits to detection performance while decreased the mAP instead such as the experimental results of boat instances.
- According to experimental results, the slot-based augmentation approach performed closed behaviour as image flipping with less increased images. Besides slot-based augmentation has large potential to generate images which is a controllable progress by customising slot match ratio and the iteration amount.
- The combination of flipping and slot-based method is relatively more robust and contributes a higher mAP.

Despite those benefits of slot-based augmentation methods, there is a notable effect on contextual features. In this chapter, augmentation methods are discussed on improving detection performance while there are many DNN architectural works aiming at extracting the features between foreground and background to enhance detection performance, such as the relation model proposed by [13]. Slot substitution crops original foreground and replace it with a new one which changes the foreground-background relational features. However, this type of changes has critical effects. On one hand, breaking the original contextual features makes it difficult for DNN models to detect objects according to their surroundings. On the other hand, breaching original contextual features potentially contributes a robust feature extraction that recognising objects fully by its features not by the environments. For example, human beings recognise the flying beer on the sky from advertisement posters.

## References

- [1] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [2] Hankyeol Lee, SEOK EON CHOI, and Changick Kim. A memory model based on the siamese network for long-term tracking. In *European Conference on Computer Vision Workshop*. Springer, 2018.
- [3] Samuel Schulter, Paul Vernaza, Wongun Choi, and Manmohan Chandraker. Deep network flow for multi-object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6951–6960, 2017.
- [4] G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.

- [5] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikit-image: image processing in python. *PeerJ*, 2:e453, 2014.
- [6] Daniel Ho, Eric Liang, Ion Stoica, Pieter Abbeel, and Xi Chen. Population based augmentation: Efficient learning of augmentation policy schedules. *arXiv preprint arXiv:1905.05393*, 2019.
- [7] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- [8] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- [9] Nikita Dvornik, Julien Mairal, and Cordelia Schmid. Modeling visual context is key to augmenting object detection datasets. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 364–380, 2018.
- [10] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection snip. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3578–3587, 2018.
- [11] Mate Kisantal, Zbigniew Wojna, Jakub Murawski, Jacek Naruniec, and Kyunghyun Cho. Augmentation for small object detection. *CoRR*, abs/1902.07296, 2019.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1904–1916, 2015.
- [13] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3588–3597, 2018.
- [14] Zhe Chen, Shaoli Huang, and Dacheng Tao. Context refinement for object detection. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [15] 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 *European conference on computer vision*, pages 740–755. Springer, 2014.
- [16] Xinlei Chen and Abhinav Gupta. An implementation of faster rcnn with study for region sampling. *arXiv preprint arXiv:1702.02138*, 2017.
- [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [19] Stefan Van Der Walt, S Chris Colbert, and Gael Varoquaux. The numpy array: a structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22, 2011.
- [20] John D Hunter. Matplotlib: A 2d graphics environment. *Computing in science & engineering*, 9(3):90, 2007.
- [21] Wes McKinney. pandas: a foundational python library for data analysis and statistics. *Python for High Performance and Scientific Computing*, 14, 2011.
- [22] 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*, pages 91–99, 2015.