
Reproducibility report of "Mind the pad"

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 Scope of Reproducibility

3 Convolution mechanism is widely adopted for a large variety of tasks like image classification or object detection.
4 Alsallakah et al. [1] demonstrate how this mechanism has some flaws caused by padding. Our aim is to reproduce the
5 following results

- 6 • Single Shot Detector blind spot on small object detection
- 7 • Single Shot Detector blind spot fix when changing padding mode from *zeros* to *reflect*
- 8 • Uneven application of padding on downsampling convolutional layers causes feature map erosion and lower
9 accuracy.

10 Once reproduced these results, we performed a series of ablation studies to understand the effect of related factors in a
11 CNN

- 12 • How does Batch Normalization interact with uneven application of padding?
- 13 • Which category between these: 1) padding modalities {*zeros*, *reflect*, *circular*, *replicate*} 2) with/without batch
14 normalization 3) with/without uneven application of padding is more shift robust on image classification?

15 Methodology

16 To reproduce paper claims we implemented all the experiments from scratch in order to have more reliable results.
17 The only external resource used in this article is a PyTorch implementation of Single Shot Detector made by NVIDIA
18 [6]. Furthermore, since the paper thesis is related not to a specific implementation but to a class of models, our
19 implementations fit inside the same category but with a different configuration. We have done so to stress paper claims
20 and to confirm their general validity.

21 To train the models for image classification, 48, we used a local Nvidia GeForce GTX 1660 with 6 GB of memory, 8
22 GB of RAM, and AMD Ryzen 5 2600X (12) @ 3.600GHz. The training holds for approximately 20 hours. The reason
23 to train such quantity of models are: 1) 4 padding modes 2) with/without Batch Normalization 3) with/without uneven
24 application of padding (i.e. input images zero/one padded) 4) 3 random seeds

25 Results

26 During our experiments, we found the same blind spots of paper authors but on a different location: for us, the blind
27 spot was close to the right border while paper claims it is at the top border. We fixed blind spot issue using *reflect*
28 padding instead of *zeros* like paper does. About uneven application of padding, we have a performance improvement
29 when comparing models with/without uneven application of padding but with different delta. In our experiments, the
30 accuracy improvement is around 0.8% while averagely of 0.4 % for the original paper. We believe that the cause is
31 a different model architecture and dataset. In this article, we adopted a simple Sequential CNN classifier on Letter
32 MNIST while the paper uses famous architectures like ResNet on Imagenet.

33 What was easy

- 34 • Obtain paper results about image classifier uneven application of padding worked at the first shot.
- 35 • Once found good implementation of SSD reproduce results on evaluation was immediate.

36 What was difficult

- 37 • Find an image classifier architecture suitable for uneven application of padding tests like the article i.e. with
38 original input size the downsampling layers don't see the right padded border while with one-pad images the
39 downsampling layers see all padded borders.
- 40 • Change padding mode of convolutional layers or disable Batch Normalization layers especially on Tensorflow
41 models.
- 42 • Understand how to plot SSD's object confidence for all zones of the image

43 Communication with original authors

44 To make this article there wasn't any contact with the original authors.

45 1 Introduction

46 Convolution algorithm is adopted in a large variety of tasks, such as image classification, object detection, and
47 generative models. However, despite its diffusion, there are a lot of hidden mechanisms we don't fully understand about
48 convolution. For instance, Alsallakah et al. [1] noted that an object detector, SSD [5], fails to recognize traffic lights on
49 consecutive frames while previously the object was identified. The major difference between frames was a vertical
50 shift of the traffic light. Starting with this observation, they demonstrate how padding can create blind spots on object
51 detection tasks. Furthermore, on image classification, they shown how downsampling layers (convolutions with stride >
52 1) don't see the right border leads to lower accuracy. They also proposed a series of constrains to avoid that phenomena,
53 called by them uneven application of padding.

54 The article focuses on reproducing the same results of Alsallakah et al. [1] through different datasets and models but in
55 the same task category. For example, instead of detecting blind spots on traffic light [2] detection with SSD [5], we
56 tested Alsallakah et al. [1] blind spot thesis using COCO dataset on images containing only small items. In the second
57 part, we trained an image classifier from scratch on Letter MNIST to verify if the uneven application of padding leads to
58 feature erosion. Furthermore, we provided a set of ablation studies to evaluate the effectiveness of Batch Normalization
59 and tested shift robustness.

60 2 Object Detection with SSD

61 2.1 Introduction

62 Single Shot Detector [5] (SSD) is an object detection model, able to identify items of different scale and size inside the
63 image in a single step. This allows the model to be easily deployable in a real time environment.

64 In our experiments, we used SSD [6] trained on COCO instead of traffic light [2], as Alsallakah et al. [1] does. For
65 evaluation we took traffic light [2] images like Figure 2 from COCO to fit in the same category of Alsallakah et al. [1],
66 i.e. small object detection.

67 2.2 Results

68 Alsallakah et al. [1] demonstrate how padding allows the model to exploit border locations causing identification
69 problems when close to the borders. Figure 3 shows very similar results of the original paper. The blind spot differs just
70 for the location: in SSD [5] trained on COCO the blind spot is close to the right border while for the original authors is
71 close to the top border.

72 We found coherent results when changing padding mode from *zeros* to *reflect*, as Figure 3 shows. We tried as well
73 padding mode *circular* but it isn't very effective to counter blind spots on SSD [5] evaluation.

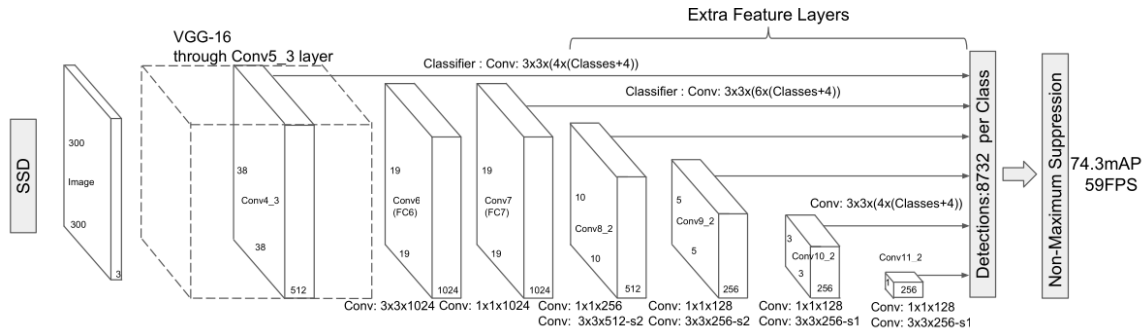


Figure 1: Single Shot Detector architecture taken from [5]. Our version has MobileNet [4] instead of VGG

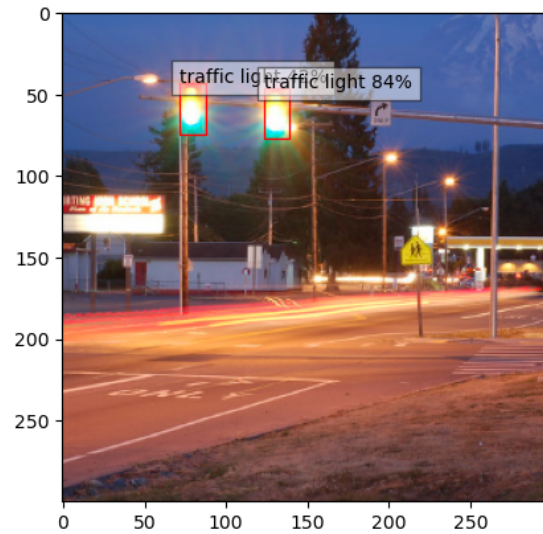


Figure 2: Sample of SSD [5] object detection with confidence above each bounding box

74 3 Spatial bias on image classifiers

75 3.1 Experiment setup

76 In order to check if Alsallakah et al. [1] thesis is correct, we recreated the same experiment condition but with different
77 models and datasets.

78 **Dataset** Dataset is Letter MNIST, a set of grayscale {28, 28} images showing handwritten alphabet letters with 26
79 different classes. Train split has 60 000 samples while test split has 10 000.

80 **Model** We trained a set of small CNN feature extractors followed by Average Pooling and two Fully Connected layers.
81 Every convolutional layer has "same" padding and downsample every two (Figure 5). Given this architecture, cardinal
82 design choices are:

- 83 • Average pooling after Convolutional layers instead Max pooling because it preserves shift equivariance
84 property and therefore, results are more reliable.
- 85 • Same padding because it's the most common in CNN. It allows the model to preserve the image size through
86 convolutional layers. This is essential for architectures like ResNet [3] or to handle variable input size.

87 We followed the same Alsallakah et al. [1] uneven application of padding (Figure 6) experiments idea. The model
88 architecture is chosen such that downsampling layers see all padded borders (no uneven application of padding) when

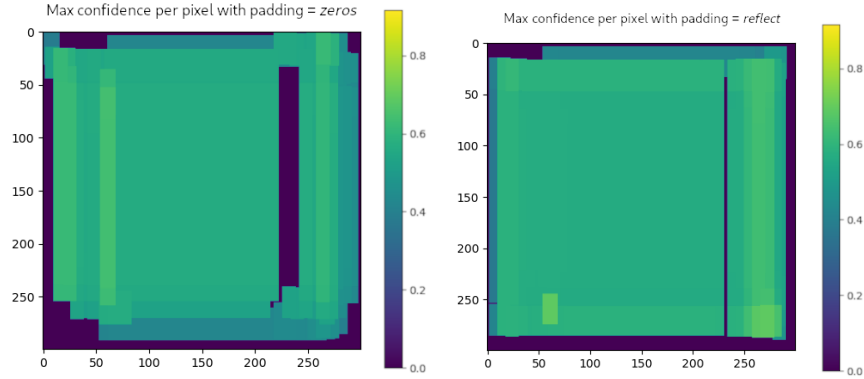


Figure 3: Example of blind spot on Figure 2 obtained by moving the left semaphore through all the image. At evaluation time SSD [5] has a blind spot close to the right border while Alsallakah et al. [1] found that on top border

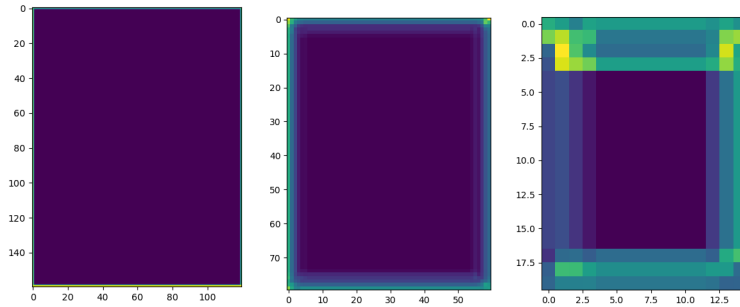


Figure 4: SSD ReLU output when feed a black image. From the left to the right, we have respectively ReLU outputs of 2nd, 9th and 19th layers. It is clear how going deeper through layers worsen border bias.

89 input size is $\{29, 29\}$ i.e. padded one. Instead with the original input size $\{28, 28\}$, downsampling layers don't see the
 90 right padded border leading to feature maps erosion.

91 **Training procedure** We have chosen the simpler training possible: model parameter optimization continues until
 92 validation loss decreases (Early stopping) using Adam optimizer with learning rate 10^{-3} .

93 3.2 Accuracy of models without uneven application of padding

94 To validate paper results, we do not aim to have the same absolute value of accuracy because model architecture, dataset,
 95 and the number of classes differ. But rather we aimed to the same relative accuracy improvement when fixing the
 96 problem of uneven application of padding.

97 As Table 1 shows, we have slightly better results to Alsallakah et al. [1] in terms of relative accuracy improvement. When
 98 uneven application of padding is missing, accuracy improves averagely of 0.8%. We believe that our improvements are
 99 slightly higher because we have done the same experiments with a simpler dataset and model architecture compared to
 100 the original paper. In their results, they showed various versions of ResNet [3] and MobileNet [4] trained on ImageNet
 101 [7].

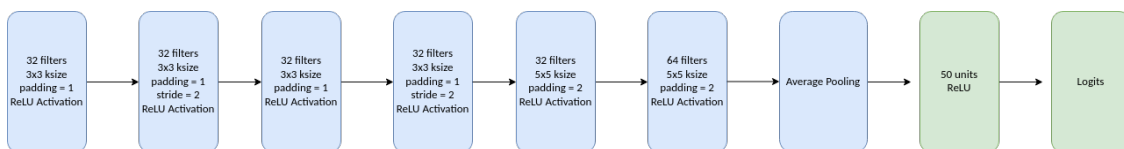


Figure 5: CNN classifier with same padding used in Section

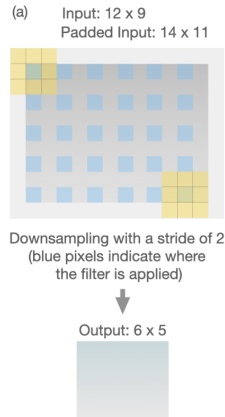


Figure 6: Example of uneven application of padding. Courtesy of Alsallakah et al. [1]

Padding mode	circular	reflect	replicate	zeros	Marginal
Input size					
{29, 29}	92.7 ± 0.2 %	93.1 ± 0.2 %	93.0 ± 0.3 %	93.0 ± 0.2 %	93.0 ± 0.2 %
{28, 28}	92.0 ± 0.3 %	92.3 ± 0.3 %	92.2 ± 0.3 %	92.1 ± 0.3 %	92.2 ± 0.3 %

Table 1: Accuracy of CNN classifier by padding mode and presence of border asymmetry in downsampling layers. Metrics results are in the form $\mu \pm \sigma$ because for each cell in the table we trained 3 identical models but with different random seed in order to prevent the influence of aleatory factors like batch sampling or parameter initialization

102 Instead, we haven't seen big accuracy variations with any type of padding mode. Only padding mode *circular* is slightly
 103 worse in terms of accuracy compared to the others, around 0.3%. We tested shift robustness of models with any padding
 104 mode, resulting in *circular* much more robust than others. Further details are in ablation studies section.

105 **Test only uneven application of padding** We check if Alsallakah et al. [1] uneven application of padding thesis
 106 holds only at evaluation time. We test models trained with uneven application of padding without it and vice-versa.

Test input size	{28, 28}	{29, 29}
Training input size		
{29, 29}	92.0 ± 0.2 %	93.0 ± 0.2 %
{28, 28}	92.1 ± 0.3 %	72.9 ± 2.4 %

Table 2: Accuracy of CNN classifier on Letter MNIST fed with border asymmetry images (size {28, 28})

107

108 Table 2 shows a performance decrease from 92.1% to 72.9% when the model trained with the uneven application of
 109 padding is evaluated without it. This result eliminates the hypothesis that just uneven application of padding leads
 110 to performance decrease. Rather, the true cause is the uneven application of padding training leads the model to
 111 learn skewed filters and therefore to worse accuracy. Furthermore, it confirms Alsallakah et al. [1] binding between
 112 asymmetrical filters and performance decrease.

Training Input size With batch norm	{28, 28}	{29, 29}	Marginal
True	14.612547	13.811748	14.212148
False	15.558550	14.280053	14.919302
Marginal	15.085549	14.045901	

Table 3: Perplexity of CNN classifiers with zero/one pad and with/without Batch Normalization. Model with Batch Normalization and one padded tends to have less perplexity at test time

113 3.3 Code to check if model has uneven application of padding

114 The following python function tells if a model with a fixed input size has uneven application of padding on downsampling
 115 layers. This implementation uses the same equations of analyzed paper Section 5. It supports *PyTorch* models and
 116 leverages on the package *Torchinfo*, which tells input and output size of each CNN’s layer.

```

import torch.nn as nn
from torchinfo import summary

def has_uneven_padding(model: nn.Module, input_size: tuple):
    summary_info = summary(model, input_size, verbose=0)
    print(summary_info)
    for layer_info in summary_info.summary_list:
        if isinstance(layer_info.module, nn.Conv2d):
            conv_module = layer_info.module
            if any(s > 1 for s in conv_module.stride): #downsampling layers
                *, h_i1, w_i1 = layer_info.input_size
                *, h_i, w_i = layer_info.output_size
                h_i_line = h_i1 + 2 * conv_module.padding[0]
                w_i_line = w_i1 + 2 * conv_module.padding[1]

                h_i_hat = conv_module.stride[0] * (h_i - 1) + conv_module.kernel_size[0]
                w_i_hat = conv_module.stride[1] * (w_i - 1) + conv_module.kernel_size[1]

                if h_i_line != h_i_hat or w_i_line != w_i_hat:
                    print('layer', layer_info.var_name)
                    print(f'{h_i_line = }', f'{h_i_hat = }', f'{w_i_line = }', f'{w_i_hat = }')
                    return True
    return False

```

118 4 Ablation studies

119 4.1 Batch Normalization impact

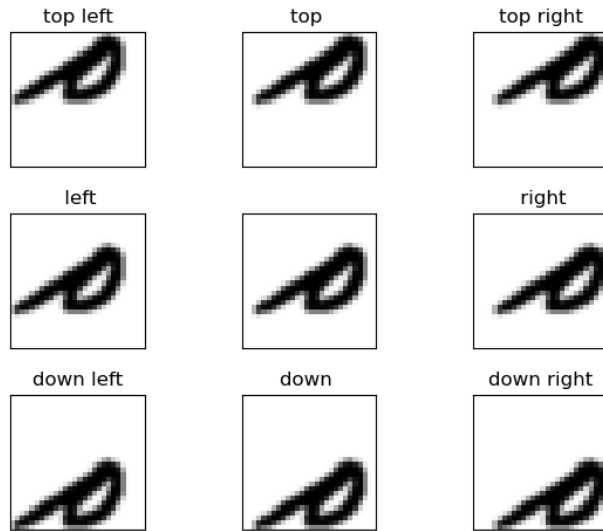
120 Once asserted coherent results with Alsallakah et al. [1], we evaluate the impact of other essential CNNs components
 121 like Batch Normalization. Using the same CNN classifier of Section 3, we added Batch Normalization after every
 122 Convolutional layer.

123 We **haven’t observed a significant difference in terms of accuracy** of the model with Batch Normalization. Still, they
 124 tend to have lower perplexity than models without Batch Normalization as we can see in Table 3.

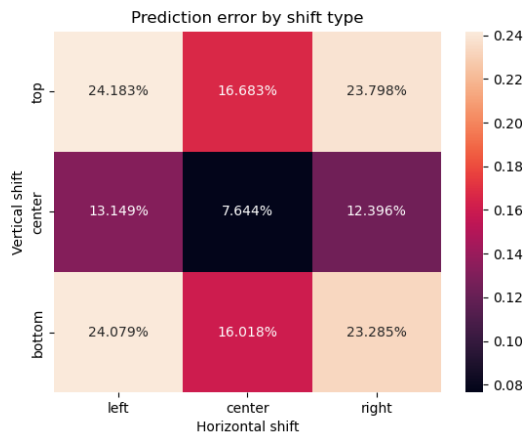
125

126 4.2 Shift robustness and prediction error

127 Following Alsallakah et al. [1] paper core principles, we analyzed *spatial bias* of image classifiers through shift
 128 robustness. With the latter concept we intend to pose the following question: the model is able to classify correctly the
 129 image even if input is shifted towards the borders?



(a) Example of shifted letter "a" where the original one is on the center. Shifts are computed algorithmically such that they push the letter at one border between $\{top, left, down, right\}$



(b) Prediction error when image is shifted toward borders. Average error is symmetric with regards to the center, except for a slight higher error when shifted to the left instead of right.

130 **Shift methodology** To generate a shifted dataset, we pushed the letter into border through any combination of
 131 these directions $\{north, center, south\} \times \{left, center, right\}$ like in Figure 7a. Therefore, we evaluated the shifted
 132 images with the trained models of Section 3 and measured how the prediction error changes when shifting the letter
 133 toward border, as Figure 7b shows.

134 **Results** After shifting in 9 directions 10 images per class we computed prediction error by models with/without
 135 Batch Normalization and with/without uneven application of padding (Figure 8). We found that models without uneven
 136 application of padding and with Batch Normalization are slightly more shift robust than other categories. Furthermore,
 137 when horizontal shift is very large, for instance with letters like *l*, all models almost random guess the label.

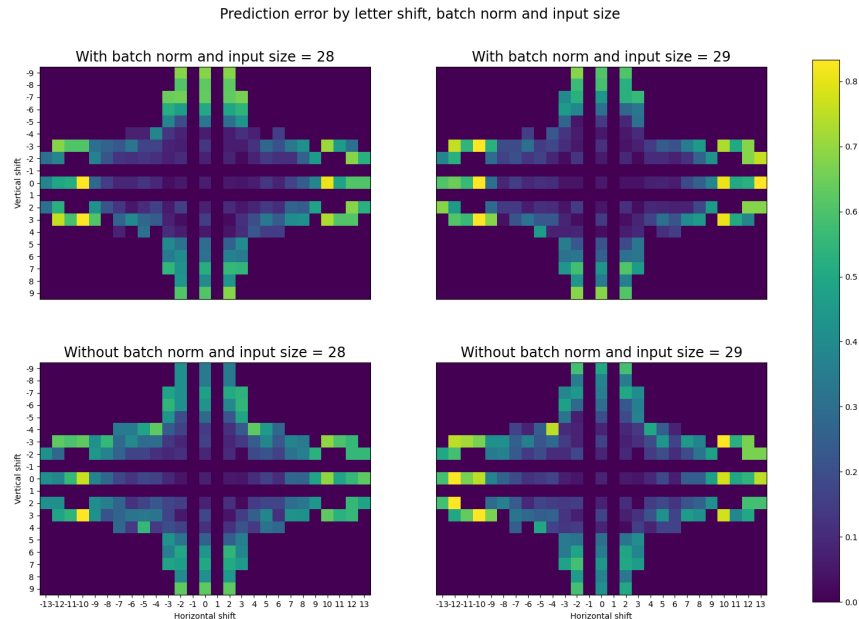


Figure 8: Prediction error by models with or without Batch Normalization, zero or one input padded. Each heatmap shows the prediction error got by shifting images toward borders as Figure 7a. Heatmap values follows a "+" pattern because letters are exclusively big vertically or horizontally. Results show that models with Batch Normalization and one pad are slightly more robust than the others.

138 Therefore we evaluated shift robustness by padding modes (Figure 9). Interestingly, models with padding mode *circular*
 139 eliminate almost all shift biases except when horizontal shift is extremely high.

140 We speculate that the primary reason for *circular* padding shift robustness is because it takes padded values from the
 141 opposite side of the feature map. By taking very distant values, the padded feature map has a higher degree of variance
 142 compared to other methods like *zeros* or *reflect* padding modes, resulting in a less "exploitable" border.

143 5 Conclusion

144 Thanks to the clarity of the paper we managed to reproduce all the main paper results with other model architectures
 145 and datasets, which makes their claims more robust. In detail, our reproduction procedures are

- 146 1. We evaluated Single Shot Detector (SSD) [5] trained on COCO dataset instead of traffic light [2] dataset
 147 obtaining compatible results. By changing SSD padding mode from *zeros* to *reflect* we fixed blind spot issue
 148 as they do.
- 149 2. We checked if image classifiers suffers from uneven application of padding using a CNN classifier built from
 150 scratch instead of using more famous architectures like ResNet [3] and MobileNet [4]. We found the same
 151 accuracy improvement when border asymmetry is removed by downsampling layers.

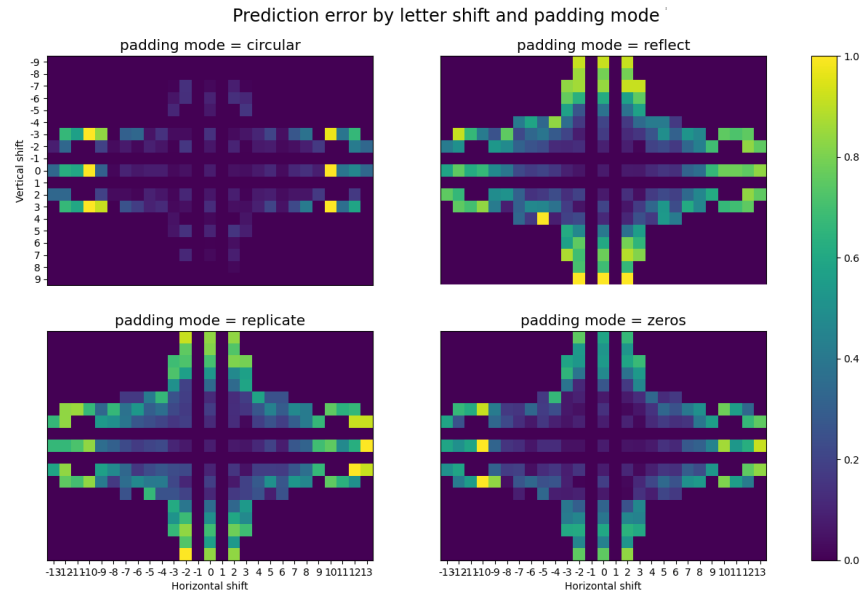


Figure 9: Prediction error by models trained with padding mode $\{zeros, replicate, reflect, circular\}$. Each plot shows prediction error got by shifting images toward borders like Figure 7a. Clearly, padding mode *circular* is much more shift robust than other modalities.

152 Following, we performed a series of ablation studies which results are

- 153 • Batch Normalization doesn't improve accuracy without uneven application of padding (Figure 6). We observed
- 154 just a slight improvement in terms of perplexity.
- 155 • We evaluated shift robustness of image classifiers by shifting the image content towards the borders through
- 156 a combination of the cardinal directions $\{north, south, east, west\}$ (Figure 7a). In this context, we found
- 157 that models with padding mode *circular* are much more shift robust than others (Figure 9) but with a slight
- 158 decrease of average accuracy (-0.3%)

159 To conclude our Reproducibility report, we propose some directions to improve actual results:

- 160 • Research new padding modalities such that they are very efficient during training and prevents blind spots on
- 161 Object detection.
- 162 • Develop new algorithms to prevent at test time skewed filters without worsening training computational cost.

163 References

- 164 [1] Bilal Alsallakh et al. "Mind the Pad—CNNs can Develop Blind Spots". In: *arXiv preprint arXiv:2010.02178*
- 165 (2020).
- 166 [2] Karsten Behrendt and Libor Novak. "A Deep Learning Approach to Traffic Lights: Detection, Tracking, and
- 167 Classification". In: *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE.
- 168 [3] Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs.CV].
- 169 [4] Andrew G. Howard et al. *MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications*.
- 170 2017. arXiv: 1704.04861 [cs.CV].
- 171 [5] Wei Liu et al. "Ssd: Single shot multibox detector". In: *European conference on computer vision*. Springer. 2016,
- 172 pp. 21–37.
- 173 [6] NVIDIA. *Object detection with Mobilenet SSD written trained on COCO*. 2020. URL: https://pytorch.org/hub/nvidia_deeplearningexamples_ssd/.
- 174
- 175 [7] Olga Russakovsky et al. "ImageNet Large Scale Visual Recognition Challenge". In: *International Journal of*
- 176 *Computer Vision (IJCV)* 115.3 (2015), pp. 211–252. DOI: 10.1007/s11263-015-0816-y.