
Reproducibility Report

RetinaFace: Single-shot Multi-level Face Localization in the Wild

Anonymous Author(s)
Affiliation
Address
email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 RetinaFace [2] is a deep learning model that detects faces in images by proposing rectangular areas (bounding boxes)
4 for every single face. Unlike the other current state-of-the-art models, this study proposes a multi-task loss calculation
5 by also computing the coordinates of 5 facial landmarks (eyes, nose, and two sides of the mouth) and 3D face mesh
6 with 1000 points concurrently. Additionally, the proposed model also adapts a cascaded structure [13] and deformable
7 convolution layers (DCL) [1]. The scope of this paper includes the whole model structure excluding DCL. Additionally,
8 The tasks implemented are limited only to face bounding box detection and landmark localization tasks, since the 3D
9 point detection database is not publicly shared.

10 **Methodology**

11 For this challenge, I implemented this model in Julia programming language, by using the Knet deep learning framework.
12 The whole model is implemented from scratch. There are official and unofficial implementations are available but these
13 codes only contain a subset of the whole model proposed in the paper. In the context module part and for constructing
14 the methods related to box proposal, these repositories are taken as examples. For training, the WIDER FACE database
15 [12] is preferred and as landmark data, custom annotations created by the original paper's authors are used. Model is
16 trained in one Tesla V100 GPU with a batch size of 10 for 60 epochs, which lasted approximately 9 days.

17 **Results**

18 The average precision (AP) metric is used for evaluation and the results are 0.093 lower in the Easy, 0.076 in the
19 Medium, and 0.129 in the Hard subsets of WIDER FACE. Possible reasons for this performance difference are discussed
20 in the *Limitations & Problems* section.

21 **What was easy**

22 Since the model only uses a small set of operations (convolution, batch normalization, unpooling, softmax, and ReLU).
23 Therefore implementing the whole model was easy except for the loss calculation part.

24 **What was difficult**

25 The selection process of which box proposals are for faces and which are for background and how to balance their
26 losses were not explained in the original paper in detail. Because of these obscurities, implementing the loss calculation
27 was difficult.

28 **Communication with original authors**

29 I contacted them to request access to the 3D face points database but learned that that data belongs to a start-up company
30 and is not publicly licensed.

31 **1 Introduction**

32 Face Detection is a crucial problem in Computer Vision. Depending on the position of the camera relative to the person
33 itself, detecting faces may bring some obstacles such as occlusions, really small-sized face captures compared to the
34 frame length, change of illumination, etc. To overcome these obstacles, some common structures are preferred in
35 the current state-of-the-art detection models. First of all, a mechanism called Feature Pyramid Network (FPN) [4] is
36 integrated for processing different-sized intermediate outputs retrieved from backbones such as ResNet [3] or VGG
37 [10]. With this approach, each intermediate output learns to focus on faces with different scales, and hence the overall
38 performance increases.

39 Secondly, two main ways are proposed for detection tasks: single-shot [5] and two-shot [8]. While the single-shot
40 method focuses on predicting the coordinates only by feeding the input to the regarding model one time, the two-shot
41 method predicts some intermediate outputs by feeding the input for the first time and more accurate bounding box
42 proposals are found by using these intermediate outputs and feeding the input to the model a second time. Although
43 two-shot models may compute more accurate results compared to single-shot models, longer processing time and extra
44 computational load directed more recent studies to design single-shot models. With the recent development and new
45 methods on single-shot object detection, current state-of-the-art models also achieved to outperform two-shot based
46 approaches [11]. Currently, three single-shot detection models RetinaFace, ProgressFace [14] and HAMBox [6] are
47 showing state-of-the-art performances in face detection. While HAMBox deals with adjusting the anchor boxes to have
48 a greater intersection of union (IOU) values with the ground truth values, ProgressFace proposes a progressive learning
49 model, which learns faces from large to small scales gradually.

50 **2 Scope of Reproducibility**

51 Current and previous state-of-the-art face detection models are bounded with using only background/face classification
52 and bounding box regression. Other tasks such as facial landmark localization or 3D face points detection are ignored,
53 which may improve the overall face detection performance if integrated into the model. Different from these studies,
54 RetinaFace aims to benefit from landmark localization and 3D face point detection tasks to improve its performance
55 further. Since 3D points data is not shared in public, this reproduction study is only limited with landmark localization
56 task added to the main detection task.

57 In this study, I am testing:

- 58 • How much the model's performance increase if landmark task is also added?
- 59 • How much improvement is seen if a context module from SSH [7] is also integrated into the baseline model?
- 60 • Does the addition of the cascaded structure improves the model performance?

61 **3 Methodology**

62 The codes of this project are shared as a supplementary material and the GitHub repository can also be shared upon
63 request.

64 **3.1 Databases**

65 This study uses only one face detection database: WIDER FACE Face Detection Database [12]. This database does not
66 include any landmark annotation by default. Therefore, 84.6k Faces in the WIDER FACE dataset are also manually
67 annotated with 5 landmarks each¹ during the original research.

68 **3.2 Architectural Details**

69 For the implementation, Julia v1.5.3 is selected as the programming language and Knet v1.4.5 is preferred as the deep
70 learning framework. Knet is a low-level framework compared to other popular deep learning frameworks such as
71 PyTorch, TensorFlow, or Keras. It only includes operation functions required during the model implementation but does
72 not provide objects for these operations. Therefore, even convolution and batch normalization (BN) layer objects are
73 constructed manually. Furthermore, Knet does not have support for multi-GPU design. Hence, the implementation is
74 only based on a single GPU usage.

¹https://www.dropbox.com/s/7j70r3eepe4r2g/retinaface_gt_v1.1.zip?dl=0

75 **3.3 Inputs & Outputs**

76 640x640 RGB images are used as the model inputs. For each of the images, confidence scores of each area proposals
 77 (one score for being background and one score for being a face), box coordinates (center coordinates, width, and height),
 78 and 5 facial landmark locations are retrieved as output.

79 **3.4 Data Augmentation**

80 As proposed in the original paper, the techniques below are used to augment the data:

- 81 • **Random Crops:** the images are cropped in square shapes randomly, then reshaped to 640x640 and the
 82 annotations are arranged concerning these crops.
- 83 • **Horizontal Flip:** With 0.5 probability, the images are flipped horizontally.
- 84 • **Color distortions:** With 0.5 probability for each, brightness, saturation, contrast, and hue distortions are
 85 applied to the input image.

86 Although data augmentation is also implemented, the different model variations are trained without any augmentation
 87 and the final fully structured model is also trained with augmentation to observe the impact of the augmentation process.

88 **3.5 Model Description**

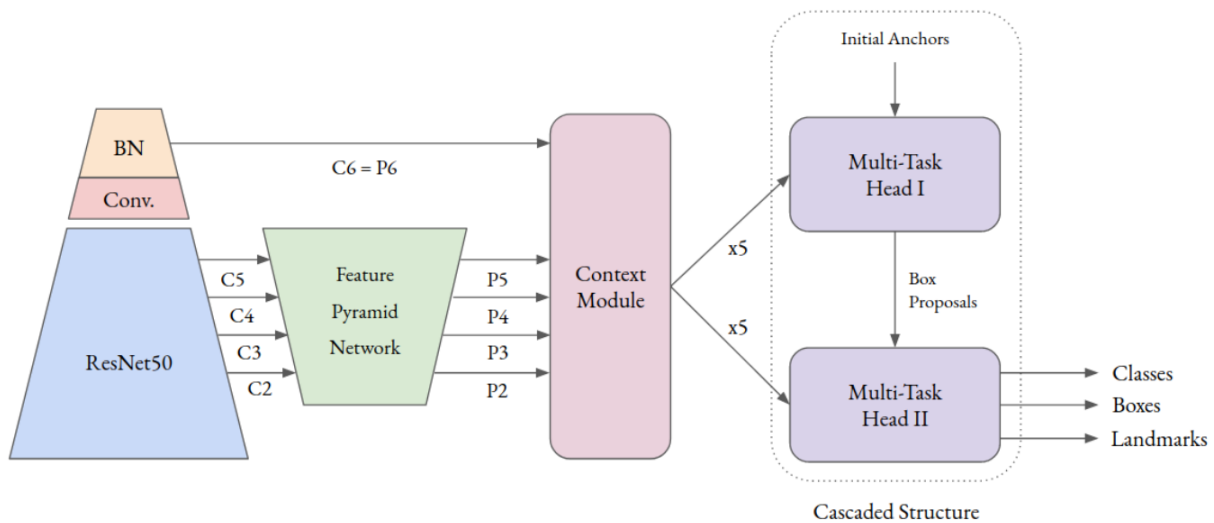


Figure 1: Structure of the Whole Model

89 **Backbone.** ResNet50 [3] is used with ImageNet pre-trained weights. From this structure, intermediate outputs of
 90 each block of convolutions are extracted. In total, there are 4 different outputs: C2, C3, C4 and C5 respectively. The
 91 network is fed with a batch of 640 x 640 RGB images in training and the dimensions below are given for each of these
 92 intermediate feature maps:

- 93 • C2: 160 x 160 x 256 x N
- 94 • C3: 80 x 80 x 512 x N
- 95 • C4: 40 x 40 x 1024 x N
- 96 • C5: 20 x 20 x 2048 x N

97 where N is the batch size. Furthermore, a 3x3 convolution + batch normalization layer (from now on, all convolution +
 98 batch normalization layer blocks will be called as ConvBn) with the stride size 2 and the filter size of 256 is also defined
 99 additionally on top of C5. The parameters of this layer are initialized with the Xavier method². With this additional

²<https://denizyuret.github.io/Knet.jl/latest/reference/#Knet.Train20.xavier>

100 layer, an extra output called C6 is created with a size of $10 \times 10 \times 256 \times N$.
 101

102 **Feature Pyramid Network.** After retrieving the outputs C2-C5, all of these values are passed on 1×1 ConvBn layers
 103 to reduce their third dimensions to 256. The new outcomes are named P2-P5. Starting from the topmost feature map
 104 (P5), an unpool (upsampling) operation is applied to equalize the first and second dimensions of the upper and lower
 105 feature maps. Then the unpooled upper layer and the lower layer are added together. Lastly, the outcome is passed to an
 106 additional ConvBn structure with a kernel size of 3×3 . C6 is excluded from all of these processes. The latest values of
 107 the intermediate feature maps can be summarized as:

- 108 • P6 = C6 (has the size: $10 \times 10 \times 256 \times N$)
- 109 • P5 = ConvBn_{1x1}(C5) (has the size: $20 \times 20 \times 256 \times N$)
- 110 • P4 = ConvBn_{3x3}(ConvBn_{1x1}(C4) + unpool(P5)) (has the size: $40 \times 40 \times 256 \times N$)
- 111 • P3 = ConvBn_{3x3}(ConvBn_{1x1}(C3) + unpool(P4)) (has the size: $80 \times 80 \times 256 \times N$)
- 112 • P2 = ConvBn_{3x3}(ConvBn_{1x1}(C2) + unpool(P3)) (has the size: $160 \times 160 \times 256 \times N$)

113 where each of the ConvBn layers is defined independently from each other. From now on, all of the P2-P6 outputs will
 114 be referred to as lateral paths.

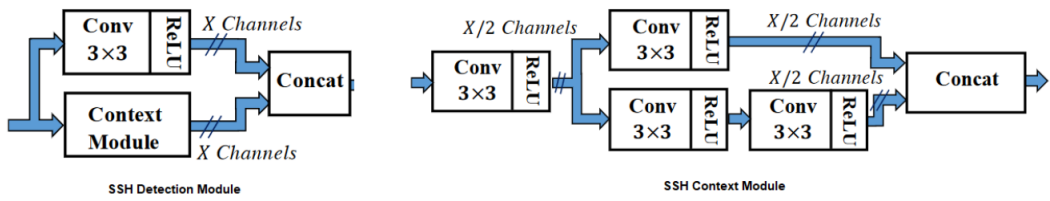


Figure 2: SSH Detection Module [7]

115 **Context Modules.** Each of the lateral paths is passed into independent context modules, which have the same model
 116 structure but different weights (in total, there are 5 context modules, one for each lateral path). The context module
 117 design is directly adopted from SSH Detection Module [7] and it can be further investigated in Figure 2. Additionally,
 118 batch normalization layers are added after each of the convolution layers in the figure. The input size is preserved by
 119 adding padding of size 1 to each convolution. All of the lateral paths are updated with this module.

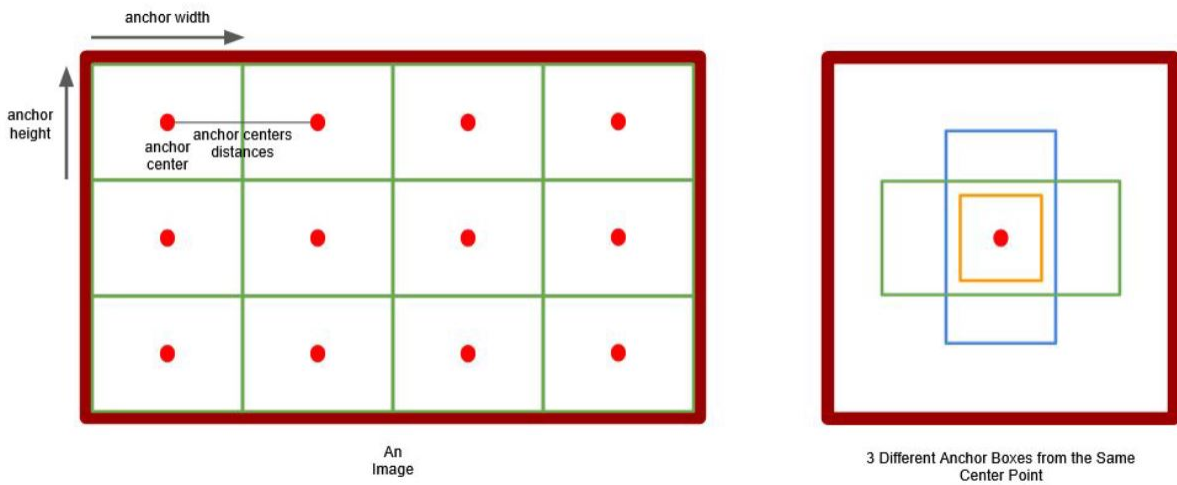


Figure 3: Anchor Box Demonstration

120 **Anchor Boxes.** Anchor boxes are predefined rectangles that divide the image into grids with different sizes. They have
 121 basically a height, width, and two center coordinates (x and y). Figure 3 demonstrates how anchor boxes are located.
 122 There is a constant distance between the centers of each neighboring anchor box pairs. Also, more than one anchor
 123 boxes can be constructed from the same center point. Each anchor box includes 4 values (center x , center y , width, and
 124 height).

125 The size of the model output is equal to the size of the combination of anchor boxes. In this specific implementation,
 126 there are 102,300 anchor boxes for a single image, the final outputs of the model are $2 \times 102,300$ for classification, $4 \times$
 127 $102,300$ for box proposal, and $10 \times 102,300$ for landmark localization tasks. By calculating the IOU values of each
 128 of the anchor boxes with each of the ground-truth bounding boxes and by selecting the maximum-IOU-value-giving
 129 anchor indices per ground truth face, the indices in the model output which are responsible for predicting the bounding
 130 box values of the ground truth objects can be determined.

131 Since faces usually share similar height and width ratios, selecting anchor boxes that have 1:1 height and width ratios
 132 would be logical. In RetinaFace, each lateral path output includes $W \times H$ many anchor box centers if the dimension of
 133 the lateral output is $W \times H \times C \times N$. Furthermore, each anchor box center is responsible for 3 different anchor boxes, all
 134 having 1:1 width and height ratios. The details of each lateral outputs are as follows:

Laterals	Anchor Center Sizes	Center Distances (in Pixels)	Edge Lengths (in Pixels)		
P6	100	64	256	322.54	406.37
P5	400	32	128	161.26	203.19
P4	1600	16	64	80.63	101.59
P3	6400	8	32	40.32	50.80
P2	25600	4	16	20.16	25.40
Total Centers	34100	Total Anchors	102300		

Table 1: Anchor Data for Each Lateral Path

135 **Multi-Task Heads.** After passing the context modules, the lateral paths are put 1×1 convolutions to convert their
 136 third dimension size to the correct output format. Assuming that $W \times H \times C \times N$ is the size of a lateral path after its
 137 corresponding context module and each anchor center is responsible for "A" many anchors (A is selected as 3 in the
 138 original paper), the final outputs retrieved from only one lateral path should have the dimension:

- 139 • Classification Output: $W \times H \times 2A \times N$
- 140 • Bounding Box Output: $W \times H \times 4A \times N$
- 141 • Landmark Output: $W \times H \times 10A \times N$

142 Multi-task heads convert the dimension C to required sizes by 1×1 convolutions. Each lateral path has its own
 143 classification, bounding box, and landmark head and after this process, each lateral path creates 3 different outputs.
 144 However, the processes after multi-task heads are common to each lateral path. Therefore, these outcomes are
 145 concatenated. To achieve this, classification outputs are reshaped to the size $2 \times (W \cdot H \cdot A) \times N$, bounding box outputs
 146 to $4 \times (W \cdot H \cdot A) \times N$ and landmark outputs to $10 \times (W \cdot H \cdot A) \times N$. Afterward, the results of each lateral path are
 147 concatenated along their second dimensions. In the end, only 3 outputs are constructed:

- 148 • Classification Final Output: $2 \times 102300 \times N$
- 149 • Bounding Box Final Output: $4 \times 102300 \times N$
- 150 • Landmark Final Output: $10 \times 102300 \times N$

151 where $(W \cdot H \cdot A)$ is the number of anchor boxes responsible for one lateral path output, 102,300 is the total number of
 152 anchor boxes.

154 **Prediction & Ground Truth Conversions.** The predicted box and landmark coordinates are designed to be scale
 155 invariant. Therefore, some transformation equations are applied to convert ground truth data to predicted value format.
 156 For any ground truth with a center coordinate (x_c, y_c) , with the lengths (w, h) and with a landmark point (x_l, y_l) :

- 157 • $x_c^p = (x_c - x_c^a)/w^a$
- 158 • $y_c^p = (y_c - y_c^a)/h^a$
- 159 • $w^p = \ln(w/w^a)$
- 160 • $h^p = \ln(h/h^a)$
- 161 • $x_l^p = (x_l - x_c^a)/w^a$
- 162 • $y_l^p = (y_l - y_c^a)/h^a$

163 where (x_c^p, y_c^p) corresponds to the box center in the prediction format, (w^p, h^p) to the box lengths in the prediction
 164 format, and (x_l^p, y_l^p) to the landmark point in the prediction format. Additionally, (x_c^a, y_c^a) is the center coordinate and
 165 (w^a, h^a) are the lengths of the corresponding anchor box.

167 **Online Hard Example Mining (OHEM) [9] & Loss Calculation.** The training process utilizes both classification
 168 and regression losses. For classification, the negative log-likelihood (NLL) is preferred and for regression, the smooth
 169 L1 loss is selected as proposed in the original paper.

$$\text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise,} \end{cases}$$

170

171 The regression losses are computed only from the positive anchors (the ones that are matched with a ground-truth
 172 object) since there are no values to regress for backgrounds. The box losses are computed by using width, height, and
 173 center coordinate values. The overall loss structure for a single selected index i is given below:

- 174 • $L_i = L_{cls}(p_i, p_i^*) + p_i^* L_{box}(t_i, t_i^*) + p_i^* L_{pts}(l_i, l_i^*)$
- 175 • $L_{cls}(p_i, p_i^*) = -\ln(1 - |p_i - p_i^*|)$
- 176 • $L_{box}(t_i, t_i^*) = \text{smooth}_{L_1}(|t_i - t_i^*|)$ for each (x_c, y_c, w, h)
- 177 • $L_{pts}(l_i, l_i^*) = \text{smooth}_{L_1}(|l_i - l_i^*|)$ for each (x, y) of 5 landmarks

178 where p_i^* is 1 if the proposal belongs to a ground-truth face and 0 otherwise, p_i means the probability of a bounding
 179 box proposal to be positive, t_i^* is a ground truth box value converted as explained in the *Prediction & Ground Truth*
 180 *Conversions* subsection, t_i is a box value prediction, l_i^* is a ground truth landmark value converted in a scale-invariant
 181 style, l_i is a landmark value prediction.

182 There is a significant difference in terms of counts of the positive anchors and negative anchors (the ones that are not
 183 assigned to a specific ground-truth object). Therefore, the OHEM method is used to balance the ratios of positive and
 184 negative anchors selected.

185 After calculating the output of the multi-head module, the IOU values are calculated between each ground truth object
 186 and each anchor box of the corresponding multi-head module. Then, the maximum IOU value is calculated for each of
 187 the anchor boxes. The anchor boxes that have an IOU value bigger than the positive threshold are selected as positive
 188 anchors. The positive threshold for the first multi-task module is set as 0.7 and for the second multi-task module as 0.5.

189 Similarly, the ones among the non-positive anchors that have a maximum IOU less than the negative threshold with any
 190 of the ground truth object are chosen as negative anchor candidates. The negative threshold is assigned as 0.3 in the first
 191 multi-task head and 0.4 in the second multi-task head. According to the OHEM rule, the number of selected negative
 192 anchors must be equal to at most 3 times the number of positive anchors. Hence among the candidates, a subset having
 193 the greatest NLL loss values is selected as negative anchors.

194

195 **Cascaded Structure.** Instead of loading the whole final proposal job to only one multi-task head structure, a cascaded
 196 model divides this job into 2 multi-task head modules. Same as before, the model retrieves lateral outputs until the
 197 context module. Then by using already-existing static anchor boxes and the first multi-task head module, it produces
 198 classification scores, bounding box, and landmark predictions. A loss from this process is also calculated. Then, instead
 199 of using the actual static anchor boxes, the model uses the bounding box predictions calculated from the first multi-task
 200 head module as anchors and applies the same process this time with the second multi-task head module. The outcomes

201 of the second multi-task head become the final outputs of the model and the sum of the losses calculated from both
202 multi-task head modules are added for computing the final loss.

203 3.6 Computational Requirements

204 The code is tested on both Windows and Linux and confirmed that it is fully functioning. To obtain the best performance,
205 Tesla V100 GPU is suggested since it can run up to 10 images per batch. The model is also tested in Tesla T4 and it is
206 seen that batch size can be set at most 4. The full model completes 6.5 epochs per day if the batch size is set as 10 and
207 discarding the cascaded structure increases this amount to 10 epochs per day.

208 4 Experiments & Results

209 4.1 Hyper-Parameters

210 As it is given in the actual paper, the SGD optimizer (momentum: 0.9, weight decay: 0.0005) is selected for training
211 purposes. Since this implementation remained limited with single GPU usage, the batch size is selected as 10, which is
212 the highest amount of image count in a batch that a Tesla V100 GPU memory supports while training.

213 The learning rate is set to 0.001 between the 1st and 27th epochs, 0.004 between 28th and 39th, 0.001 again between
214 40th and 49th, and 0.0001 between 50th and 60th. In total, each of the model variations is trained for 60 epochs and the
215 individual checkpoints are chosen as final, where the highest scores are achieved. The learning rate of 0.01 is not used
216 although it is preferred in the original paper, because the batch size is 3 times lower in this implementation compared to
217 the original training batch size, and selecting high learning rates may cause unstable updates.

218 4.2 Different Model Variations

219 To test the hypotheses mentioned in the *Scope of Reproducibility* section, different variations of the model are trained
220 separately:

- 221 • **Baseline:** ResNet50 + FPN + Landmark Localization Task
- 222 • **Context Module without Landmark:** ResNet50 + FPN + Context Module
- 223 • **Context Module with Landmark:** ResNet50 + FPN + Context Module + Landmark Localization Task
- 224 • **Full Model:** ResNet50 + FPN + Context Module + Landmark Localization Task + Cascaded Structure

225 4.3 Evaluation of AP in WIDER FACE Validation Data

226 In the actual paper, the evaluation also included the performance of the landmark localization task but in this reproduction
227 study, the evaluation scope is limited only to bounding box prediction performance and the average precision (AP)
228 metric used for this purpose. AP is calculated by taking the IOU threshold as 0.5 and iterating through 1000 steps of
229 confidence levels between 0 and 1. The models are evaluated only with WIDER FACE validation data and this data is
230 separated into 3 groups (Easy, Medium, and Hard) concerning their difficulty. The confidence threshold is set to 0.02 to
231 decrease the total computation time, then the top 5000 predictions are selected among the candidates and lastly, the
232 non-maximum suppression (NMS) method is applied with a threshold of 0.4 to eliminate redundant area proposals.

233 In table 2, the results of the original RetinaFace model, other state-of-the-art models, and my different model variations
234 are provided. While the HAMBox model achieves the highest performance in all of the 3 subsets of WIDER FACE
235 validation data, RetinaFace performs close to HAMBox.

236 The best-performing model results in 0.093 lower AP value in the Easy subset, 0.076 lower in the Medium, and 0.129
237 lower in the Hard subset compared to the original paper results. As mentioned in the *Scope of Reproducibility* section,
238 the model performs better when landmark localization task is included or the Context Module is also added. However,
239 adding a cascaded structure causes a performance drop in contrast to our expectations. Although the reason for the
240 negative effect of the cascaded structure is not clear, the possible main reasons for the performance difference between
241 the original paper and this implementation are discussed in the following *Limitations & Problems* section.

Model	WIDER FACE Easy	WIDER FACE Medium	WIDER FACE Hard
HAMBox [6] Baseline	0.943	0.931	0.894
HAMBox Final	0.970	0.964	0.933
ProgressFace [14]	0.968	0.962	0.918
RetinaFace Baseline	0.958	0.952	0.899
+ Context Module with DCL	0.961	0.956	0.903
+ Cascade	0.962	0.957	0.906
+ Landmark Loss	<u>0.966</u>	<u>0.959</u>	<u>0.912</u>
Baseline	0.842	0.864	0.769
Context Module without Landmark	0.865	0.878	0.767
Context Module with Landmark	0.873	0.883	0.783
Full Model	0.842	0.854	0.752

Table 2: WIDER FACE Validation Data AP Scores

242 4.4 Some Example Visual Results

243 These results are retrieved from the best performing model variation. For the prediction, the NMS threshold is set to 0.2
 244 and the confidence threshold is set to 0.5. Images are taken from WIDER FACE validation data. If the faces are not too
 245 small, then the model mostly detects the faces (Figures 4 and 5) even when there is a bad lighting (Figure 6) or the faces
 246 from a slightly different domain (a drawing in the case of Figure 7). On the other hand, if the faces are too small or the
 247 resolution is not very clear (Figures 8, 9, and 10), then the model may miss the faces in the picture.

248



Figure 4: Good Example Result I



Figure 5: Good Example Result II



Figure 6: Good Example Result III: Bad Lighting



Figure 7: Good Example Result IV: Drawing Domain



Figure 8: Bad Example Result I: Small Faces, No Faces Detected



Figure 9: Bad Example Result II: Not a Good Resolution



Figure 10: Bad Example Result III: Small Faces

249 5 Limitations & Problems

250 In this section, some limitations of the framework and the main reasons for the difference between the actual paper's
251 results and this implementation's results will be discussed.

252 **Lack of Deformable Convolution Layers.** The actual paper uses Deformable Convolutions instead of regular convo-
253 lutions for the Context Module structure. This special type of convolution achieves to "learn the offsets from the target
254 tasks, without additional supervision" [1]. Therefore, it increases the overall performance of the model. However, this
255 type of layer structure is not available in Julia or Knet, and implementing this structure requires a significant amount of
256 additional work. Thus, this improvement is left as future work.

257 **Lower Batch Size.** The original RetinaFace model uses 4 V100 GPUs and in total 32 images as a batch. However,
258 Knet does not fully support training models on multiple GPUs and the maximum number of images a V100 GPU can
259 take are 10 images. Hence, the model is limited to complete its training with 3 times lower batch size. Having a lower
260 batch size also causes more unstable updates on the model. Therefore, it becomes harder to find the optimal checkpoint
261 and the chance of finishing the training with sub-optimal local minimum increases.

262 **Knet Defaults During Learning Rate Change.** During training, each of the trainable parameters stores an additional
263 optimizer field, where the optimizer name and its specific parameters are stored. With the current Knet configurations
264 once this field is set, calling another optimizer with different parameter settings does not change the optimizer field of
265 the parameters unlike the other frameworks such as PyTorch or Tensorflow. Hence, the model continues to be trained
266 with the initial optimizer and learning rate until the end. To change the learning rate or the optimizer, each of the
267 optimizer fields of each trainable parameter has to be set to "*nothing*". I realized this problem in the last couple of days
268 of the first submission and therefore, I could only submit some results that are significantly lower than the original
269 paper.

270 **GPU-CPU Data Transfer.** While recording the gradient flow of the training batch, Knet also forbids to slice the final
271 outputs and compute a loss from these sliced sub-parts when the data is on GPU. Therefore, a constant data transfer
272 between the GPU and CPU takes place in this implementation during the loss calculation. This deficiency also increases
273 the overall run-time.

274 6 Conclusion

275 In this paper, the structure of the RetinaFace model is analyzed, the implementation process is explained, the configura-
276 tions and different experiments are shared and the results are discussed. The scope of the reproducibility is defined
277 as testing if the model's performance increases when landmark localization task, context module, and/or cascaded
278 structure are included in the model's structure.

279 According to the AP results retrieved for different model variations, it is shown that adding context module and landmark
280 task to the model increases the performance. However, including the cascaded structure resulted in a decrease in the
281 overall performance. The evaluation score difference between the original paper and this implementation is between
282 0.076 and 0.129 but the lower batch size and lack of deformable convolutions are possible reasons for this performance
283 drop. Overall, the model mostly detects faces when they are not extremely small in the given image.

284 7 Discussion

285 7.1 What Was Easy

286 As stated in the summary page and the model description, RetinaFace only uses a couple of layer and activation
287 structures. Excluding the deformable convolution layers used in the original model, it is easy to implement the whole
288 model structure (except the loss calculation) even with a low-level deep learning framework like Knet.

289 7.2 What Was Difficult

290 Implementing the loss calculation structure was the hardest task because of the lack of information in the original
291 paper. No detailed explanation was given to select the positive and negative anchors, OHEM technique was
292 stated to balance the negative and positive anchor selection but no extra hyper-parameters and instructions to use
293 the OHEM method were provided. To overcome this issue, official and unofficial implementation codes and the
294 issues opened in these repositories are checked. Additionally, the OHEM paper is read and blogs on that issue are visited.
295

296 Implementing the AP metric in Julia was also a challenging task because it is a complicated metric and for this specific
297 case, 3 different AP evaluations have to be made for each of the subsets of WIDER FACE validation data: Easy,
298 Medium, and Hard. These subsets are not separated in terms of image files but are created by selecting a subset of faces
299 for every single image. To avoid any false evaluation results, I used a python repository³ which includes all of the AP
300 evaluation processes for every single subset of the validation data.

301 7.3 Communication with Original Authors

302 I only communicated with the authors to request the 3D facial points data to extend my model with the 3D point
303 prediction task. However, they indicated that the data is not licensed public. Therefore, I excluded the 3D point detection
304 task from the model structure. Other than this, the model was explained mostly successfully in the original paper and it
305 was enough for me to implement most of the parts.

306 References

- 307 [1] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional
308 networks, 2017.
- 309 [2] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot
310 multi-level face localisation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
311 Pattern Recognition (CVPR)*, June 2020.
- 312 [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- 313 [4] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection.
314 In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 936–944, 2017.
- 315 [5] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C.
316 Berg. Ssd: Single shot multibox detector. *Lecture Notes in Computer Science*, page 21–37, 2016.
- 317 [6] Yang Liu, Xu Tang, Xiang Wu, Junyu Han, Jingtuo Liu, and Errui Ding. Hambox: Delving into online high-quality
318 anchors mining for detecting outer faces. *CoRR*, abs/1912.09231, 2019.
- 319 [7] Mahyar Najibi, Pouya Samangouei, Rama Chellappa, and Larry Davis. Ssh: Single stage headless face detector,
320 2017.
- 321 [8] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with
322 region proposal networks, 2016.
- 323 [9] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard
324 example mining, 2016.
- 325 [10] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition,
326 2015.
- 327 [11] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-yolov4: Scaling cross stage partial
328 network, 2020.
- 329 [12] Shuo Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Wider face: A face detection benchmark. In *IEEE
330 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 331 [13] Hongkai Zhang, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Cascade retinanet: Maintaining
332 consistency for single-stage object detection. *CoRR*, abs/1907.06881, 2019.
- 333 [14] Jiashu Zhu, Dong Li, Tiantian Han, Lu Tian, and Yi Shan. Progressface: Scale-aware progressive learning for
334 face detection. In *ECCV*, 2020.

³<https://github.com/wondervictor/WiderFace-Evaluation>