

FDDB-360: Face Detection in 360-degree Fisheye Images

Jianlin Fu, Saeed Ranjbar Alvar, Ivan V. Bajic, and Rodney G. Vaughan

School of Engineering Science, Simon Fraser University

Burnaby, BC, V5A 1S6, Canada

{*jfa49, saeedr, ibajic, rodney_vaughan*}@sfu.ca

Abstract

360° cameras offer the possibility to cover a large area, for example an entire room, without using multiple distributed vision sensors. However, geometric distortions introduced by their lenses make computer vision problems more challenging. In this paper we address face detection in 360° fisheye images. We show how a face detector trained on regular images can be re-trained for this purpose, and we also provide a 360° fisheye-like version of the popular FDDB face detection dataset, which we call FDDB-360.

Keywords – face detection, 360° images, deep learning, FDDB-360 dataset

1. Introduction

Face detection is a poster problem of computer vision, with applications in surveillance, security, biometrics, human-computer interaction, and other areas. With the recent advances in deep learning, the problem of face detection in conventional 2D images has largely been solved. For example, on the well-known Face Detection Data Set and Benchmark (FDDB) [8], recent models such as [10, 14] reach near 100% true positive rate with very few false positives. However, face detection in other signal domains, such as 360° images, point clouds [11], or compressed bitstreams [1, 2], has been less explored.

Our focus in this paper is on face detection in 360° images. One way to accomplish face detection in this case would be to project the 360° image onto a set of 2D images, and then employ one of the conventional face detection models on these 2D images. An approach similar to this, where an equirectangular panorama image (a post-processed form of a 360° image) is used for 2D projection and subsequent object detection, has been presented in [13]. However, the focus of [13] is generic object detection rather



Figure 1: Sample image from our FDDB-360 dataset.

than face detection. Another recent work [9] proposes to generate a spherical image and adjust convolution kernels and pooling operators to work in spherical coordinates, so that Convolutional Neural Network (CNN)-based detectors could operate directly on spherical coordinates.

Popular dual-lens 360° cameras (e.g., Ricoh Theta, Samsung Gear 360, Insta 360, etc.) have two hemispherical (“fisheye”) lenses, each projecting a 180° view onto an imaging sensor. The images read directly from the sensors are circular in appearance, and have strong barrel distortion towards the perimeter of the circle. Both the equirectangular panorama used in [13] and the spherical-coordinate image used in [9] are obtained from these fisheye images via post-processing. Clearly, it would be much more efficient to perform face/object detection directly on fisheye images; this would enable the detector to operate closer to the sensor and circumvent unnecessary processing steps. However, training face/object detectors on fisheye images is challenging because of the lack of annotated datasets of fisheye

This work was supported in part by the NSERC Strategic Project Grant STPGP 494209

images. We address this challenge by creating a fisheye-looking version of FDDB (which we call FDDB-360) and training a face detector on it. A sample image from FDDB-360 is shown in Fig. 1.

The paper is organized as follows. Creation of FDDB-360 from the original FDDB is described in Section 2. Experiments on FDDB-360 are described in Section 3, followed by conclusions in Section 4.

2. Creating FDDB-360

The purpose of FDDB-360 is to enable the training of models for detecting faces in fisheye images. To this end, we used the annotated images from the original FDDB dataset [8] and from them created a number of images that have the appearance of fisheye images. The face locations in the original FDDB are specified by ellipses, but we first converted them to rectangles to enable easier processing.

Fisheye images have least geometric distortions near the center. As we move towards the perimeter, the degree of distortions increases. Face detectors trained on conventional 2D images would likely be able to detect faces near the center of a fisheye image, but detection will become more challenging away from the center. For this reason, we wanted to create fisheye-looking images that would have sufficient number of faces away from the center, so that the detector can learn the appearance of distorted faces. The steps taken to create the new dataset are described in the following subsections.

2.1 Image extrapolation

A number of images in the original FDDB have faces near the center. To facilitate sampling the images with patches where the face locations could be arbitrary, we widened all images by 40%, by extending it 20% on both left and right as shown in Fig. 2(a). This requires image extrapolation, which is inherently difficult due to the absence of boundary conditions on three sides of the extensions.

To overcome this challenge, we used the strategy illustrated in Fig. 2(b). We created a copy of the image on the right side of the original, spaced away by 40% of the original width, and then interpolated between the two copies of the image. Though still challenging, this is an easier problem than extrapolation in Fig. 2(a) because of a larger area where boundary conditions exist. After interpolation, the interpolated area is split in half, and the right half of it (shown in red in Fig. 2) is moved to the left side of the image, to complete the extrapolated image.

Interpolation is carried out using the inpainting algorithm from [3], which is an extension of the well-known Criminisi *et al.* method [5]. While performing inpainting,



Figure 2: Image extrapolation: (a) Final extrapolated image. (b) Extrapolation carried out as interpolation between two copies of the image.

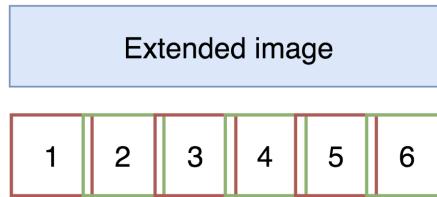


Figure 3: Six square patches, evenly spaced along the width, are extracted from an extended image.

we exclude face and skin regions from being used for inpainting. This is because, if the inpainted region ends up with patches containing partial human face, it might confuse the model during training, whether or not these partial faces are actually annotated as faces. To avoid this, while performing inpainting, we increase the cost [3, 5] of a patch for any patch overlapping with an annotated face, or where skin color is detected [4].

The approach mentioned above was applied to the majority of images in FDDB. However, about 34% of FDDB images have width-to-height ratio of less than 3:4. For such images, even 40% width extension still gives a relatively narrow image. In these cases, we did not rearrange the extended image as shown in Fig. 2(a), but left it in the format shown in Fig. 2(b), with a copy of the original image on the right side.

2.2 Fisheye-like distortion

Each extended image is sampled using square patches, which are evenly distributed along its width as shown in Fig. 3. In total, six square patches are extracted from each extended image. Subsequently, fisheye-like distortion is applied to each square patch.

Fisheye distortion models usually involve intrinsic camera parameters¹ and various lens distortion parameters [6]. To avoid making distortions camera- and lens-specific, we adopted a simpler approach.

Consider a square patch extracted from an extended image, as shown in Fig. 3. We first map this square patch to a

¹https://docs.opencv.org/3.4/db/d58/group_calib3d_fisheye.html



Figure 4: Square patch (**left**) converted to a circular patch with fisheye-like distortion (**right**).

circular patch. Let (x, y) be the normalized coordinates of the square patch, such that the patch center is $(0, 0)$ and the four corners have coordinates $(\pm 1, \pm 1)$. The square patch is mapped to a circular patch using the following coordinate mapping:²

$$(x', y') = \left(x\sqrt{1 - \frac{y^2}{2}}, y\sqrt{1 - \frac{x^2}{2}} \right). \quad (1)$$

Such mapping introduces radial distortion, where straight lines get bent towards the perimeter of the circle. To add barrel distortion, which manifests itself as “squeezing” towards the perimeter, we further scale the coordinates by a factor that decreases towards the perimeter. Specifically, the mapping is

$$(x'', y'') = \left(x'e^{-r^2/4}, y'e^{-r^2/4} \right), \quad (2)$$

where $r = \sqrt{(x')^2 + (y')^2}$ is the radial distance from the center of the patch. This form of exponential squeezing was chosen empirically to visually approximate the appearance of fisheye images. An example of a square patch converted to a circular patch with fisheye-like distortion is shown in Fig. 4.

2.3 Annotation

Once square patches are converted to circular patches, face locations have to be appropriately mapped to the new coordinates. As mentioned earlier, the location of each face in the original FDDB is specified by an ellipse, but we converted those locations to rectangles to simplify further mapping to the circular patch.

We selected eight points from the bounding box for each face (four corners and four edge midpoints), as illustrated in red in the left part of Fig. 4. If a part of the bounding box fell out of the boundaries of the square patch, we trimmed

²<https://www.xarg.org/2017/07/how-to-map-a-square-to-a-circle/>

its coordinates to coincide with the patch boundary. Once the mapping (1)-(2) is applied to these eight points, they are mapped to the squeezed circular coordinates as shown in the left part of Fig. 4. While it is possible to store these eight points (or more, if higher precision is needed) as the bounding polygon for the face in circular coordinates, we decided to simplify the annotations and again use bounding boxes. Hence, we selected the minimum bounding rectangle of the polygon as the annotation for the face location. This is illustrated by a green rectangle in the right part of Fig. 4.

When extracting square patches from the extended image, it is possible that faces are cropped and that only a part of the original face falls into the square patch. If the overlap between the original bounding box and the part that is inside the square patch is over 50%, we kept that annotation and mapped it to the circular patch as explained above. Otherwise, we treated the face as incomplete and did not convert the corresponding annotation to the circular patch.

2.4. Further details of FDDB-360

After applying the procedures described above, we ended up with 17,052 fisheye-looking images and a total of 26,640 annotated faces. Face locations are provided as the bounding box parameters (x, y, w, h) , where (x, y) is the location of the top-left corner and (w, h) are the width and height of the bounding box, respectively.

Note that the bounding box need not be fully contained within the circular patch, especially for faces that are near the perimeter, as illustrated in the right part of Fig. 4. From these coordinates, one can find the intersection of the bounding box and the circle if a more precise localization of the face is required.

The FDDB-360 dataset will be made available online at <http://www.sfu.ca/~ibajic/#data>

3. Experiments

In this section we illustrate the benefits of FDDB-360 by showing that one can re-train an existing face detection model to better detect faces in fisheye-looking images. Our baseline face detector is the well-known TinyFace [7], specifically their hybrid-resolution (HR) model, which was trained on the Wider face dataset [12]. We used their pre-trained model and tested it on FDDB-360 shown. The precision-recall curve of the HR model is shown as the green line in Fig. 5, while its ROC curve is shown as the green line in Fig. 6.

Although TinyFace HR is one of the best face detectors on 2D images according to the results on the original FDDB³, its performance on FDDB-360 is not particularly

³<http://vis-www.cs.umass.edu/fddb/results.html>

Training	Test/Validation
Folds 3-10	Folds 1-2
Folds 1-2,5-10	Folds 3-4
Folds 1-4,7-10	Folds 5-6
Folds 1-6,9-10	Folds 7-8
Folds 1-8	Folds 9-10

Table 1: Training-test split for 5-fold cross-validation

impressive. This is understandable, as the model has not been trained on images that involve fisheye-like distortions. Examination of the model’s predictions shows that indeed, near the center of the circular patches, the detection performance is quite good, but it degrades towards the perimeter where the distortions are stronger.

To show the improvement that can be obtained by transfer learning and re-training, we kept the original ten folds from the FDDB dataset. Experiments were run using 5-fold cross-validation obtained by merging pairs of the original folds as shown in Table 1. In each experimental run, we start with the pre-trained model whose weights were obtained on the Wider face dataset [12]. We then further train the model using the software provided by [7], with a learning rate of 10^{-4} on one training set shown in Table 1 and test on the corresponding test/validation set.

To account for various orientations and scales of faces that may be found in real fisheye images, we performed data augmentation during training. This included horizontal flipping and re-scaling already implemented in the software provided by [7], as well as newly introduced random rotation by 90° , 180° and 270° . We considered using other rotation angles as well. However, FDDB-360 contains annotations in the form of axis-aligned rectangles (top left x, top left y, width, height) specifying the location of the face. Rotation by angles other than 90° , 180° or 270° would cause an increase of the resulting axis-aligned enclosing rectangle, and we felt the ground-truth rectangles were already large enough (see Fig. 4) and should not be increased further. One possibility to increase the set of viable angles in augmentation in the future would be to switch to a polygonal representation of face locations.

We refer to the resulting, re-trained model, as HR-360. Red lines in Figs. 5-6 represent the average of five performance curves of HR-360 from the 5-fold cross-validation.

As seen in the figures, re-training with transfer learning from the original HR results in significant performance improvement. On the precision-recall results (Fig. 5), HR achieves the area under the curve (AUC) of 0.873, whereas HR-360 achieves 0.960, as indicated in the figure legend. On the true positives (TP) vs. false positives (FP) test (Fig. 6), HR-360 achieves around 0.85 TP rate with 200 false positives, while the original HR achieves around 0.71

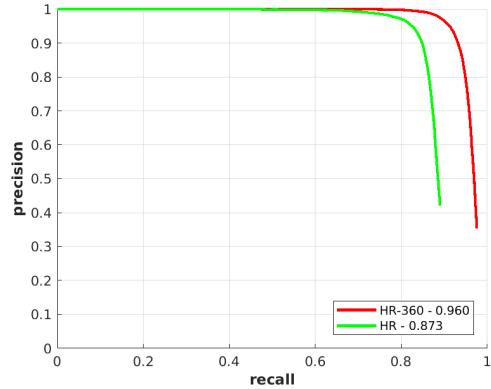


Figure 5: Precision-recall curves on FDDB-360 for HR and HR-360.

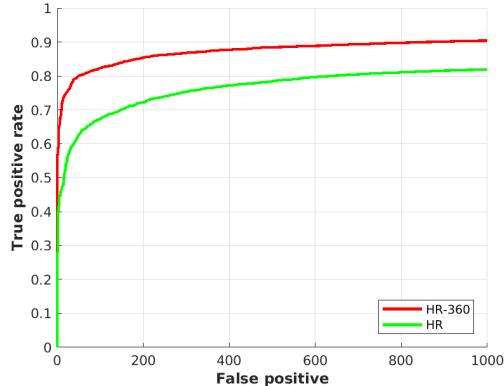


Figure 6: True positive rate versus the number of false positives for HR and HR-360 on FDDB-360.

TP rate for the same number of false positives. The difference of about 0.1 TP rate persists for higher number of false positives.

While HR-360 clearly outperforms HR on FDDB-360 overall, we did find some cases where a face was detectable by HR but not by HR-360. Apparently, transfer learning involves some “forgetting” as well, and while the model gets new capabilities during transfer learning, some of its old capabilities may disappear. One way around it could be to randomly insert the data that the model was originally trained on (in our case, the Wider face dataset [12]) into the re-training process. However, in keeping with the FDDB-style evaluation, we did not do that, and only used the data from the FDDB-360 folds (Table 1) for re-training.

To visualize how accurate are HR and HR-360 depending on the location of the target face, we test both models on the entire FDDB-360 dataset and record the locations of all false negatives (FN) - the faces that were missed. This is one way to find how a model can be improved. When a face

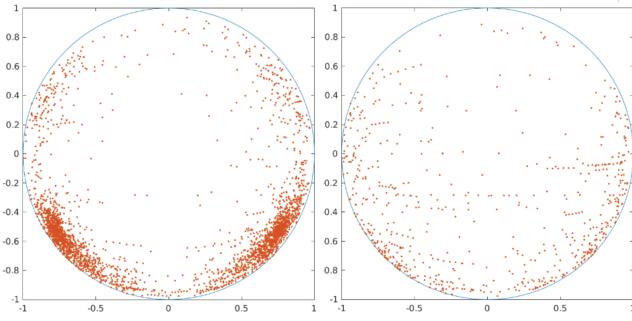


Figure 7: False negative distribution for HR (**left**) and HR-360 (**right**).

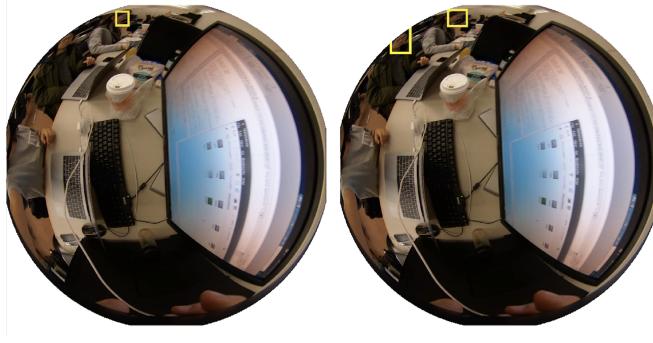


Figure 8: Face detection results by HR (**left**) and HR-360 (**right**) on a real fisheye image with several faces along the perimeter. Detected faces are indicated by yellow rectangles.

is missed, we record the location of the center of its ground-truth rectangle and normalize it in such a way that the radius of the circular image is 1 (i.e., the image becomes a unit circle). In cases where the center of the bounding rectangle is outside the unit circle, the intersection point of the circle and the line connecting the center points of the bounding rectangle and the circle is used to represent the FN location. The scatter plot of FN points for HR and HR-360 are shown in Fig. 7. The left graph shows FN’s of HR and it is apparent that HR misses many faces near the perimeter of the image, as we expected due to geometric distortions. Meanwhile, the FN’s of HR-360 (Fig. 7 right) are more evenly distributed across the unit circle, indicating that the model has learned the corresponding geometric distortions.

Finally, we show an example of detection performance on a real fisheye image, rather than an image from FDDB-360. Fig. 8 shows an image obtained by a Ricoh Theta V camera with several faces along the perimeter. Image on the left shows the results of HR, where one face was detected, as indicated by the yellow rectangle. Image on the right shows the result of HR-360, which manages to find two faces. There is one more person in the scene (in the left part of the image), whose face is so much out of view that both detectors fail to find it.

4. Conclusions

In this paper, we described the creation of FDDB-360, a dataset for face detection in fisheye images. The dataset was created from the well-known FDDB dataset by sampling patches from its images and applying fisheye-looking distortion to them, while mapping the face annotations to the new coordinate system. We also showed that re-training (using transfer learning) an existing face detector on FDDB-360 is able to significantly improve its face detection performance on this kind of images.

References

- [1] S. R. Alvar, H. Choi, and I. V. Bajić. Can you find a face in a HEVC bitstream? In *Proc. IEEE ICASSP’18*, pages 1288–1292, Apr. 2018.
- [2] S. R. Alvar, H. Choi, and I. V. Bajić. Can you tell a face from a HEVC bitstream? In *Proc. IEEE MIPR’18*, pages 257–261, Apr. 2018.
- [3] I. V. Bajić. Image inpainting with data-adaptive sparsity. In *Proc. Sinteza 2014*, pages 835–840, Belgrade, Serbia, 2014.
- [4] E. Buza, A. Akagic, and S. Omanovic. Skin detection based on image color segmentation with histogram and k-means clustering. In *Proc. Int. Conf. Electrical and Electron. Eng. (ELECO)*, 2017.
- [5] A. Criminisi, P. Perez, and K. Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Trans. Image Processing*, 13(9):1200–1212, Sep. 2004.
- [6] P. Drap and J. Lefèvre. An exact formula for calculating inverse radial lens distortions. *Sensors*, 66(6), 2016.
- [7] P. Hu and D. Ramanan. Finding tiny faces. In *Proc. IEEE CVPR’17*, pages 1522–1530, 2017.
- [8] V. Jain and E. Learned-Miller. FDDB: a benchmark for face detection in unconstrained settings. Technical Report UM-CS-2010-009, University of Massachusetts, Amherst, 2010.
- [9] Y. K. Lee, J. Jeong, J. S. Yun, C. W. June, and K.-J. Yoon. SpherePHD: Applying CNNs on a Spherical PolyHeDron representation of 360 degree images. arXiv:1811.08196, Nov. 2018.
- [10] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang. DSFD: dual shot face detector. arXiv:1810.10220, Nov. 2018.
- [11] A. R. E. Sayed, A. E. Chakik, H. Alabboud, and A. Yassine. Efficient 3D point clouds classification for face detection using linear programming and data mining. *The Imaging Science Journal*, 66(1):23–37, 2018.
- [12] S. Yang, P. Luo, C. C. Loy, and X. Tang. Wider face: A face detection benchmark. In *Proc. IEEE CVPR’16*, pages 5525–5533, 2016.
- [13] W. Yang, Y. Qian, F. Cricri, L. Fan, and J. Kamarainen. Object detection in equirectangular panorama. arXiv:1805.08009, May 2018.
- [14] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li. S³FD: single shot scale-invariant face detector. In *Proc. IEEE ICCV’17*, pages 192–201, 2017.