# LIVERPOOL JOHN MOORES UNIVERSITY

## LJMU Research Online

Khan, W, Hussain, A, Kuru, K and Al-askar, H

 Pupil Localisation and Eye Centre Estimation using Machine Learning and Computer Vision

http://researchonline.ljmu.ac.uk/id/eprint/13258/

Article

For more information please contact researchonline@ljmu.ac.uk

1 *Type of the Paper (Article)*

2 # Pupil Localisation and Eye Centre Estimation
3 using Machine Learning and Computer Vision

4

5 **Wasiq Khan [1,*], Abir Hussain [2], Kaya Kuru [3], and Haya Al-askar [4,]**

6
7 [1] Liverpool John Moores University, Liverpool, UK; W.khan@ljmu.ac.uk
8 [2] Liverpool John Moores University, Liverpool, UK; A.Hussain@ljmu.ac.uk
9 [3] University of Central Lancashire, Preston, UK; KKuru@uclan.ac.uk
10 [4] Prince Sattam Bin Abdulaziz University, Saudi Arabia; H.Alaskar@psau.edu.sa
11
12 **\*** Correspondence: W.Khan@ljmu.ac.uk;

13 **Abstract:** Various methods have been used to estimate the pupil location within an image or a
14 real-time video frame in many fields. However, these methods lack the performance specifically in
15 low-resolution images and varying background conditions. We propose a coarse-to-fine pupil
16 localisation method using a composite of machine learning and image processing algorithms. First,
17 a pre-trained model is employed for the facial landmark identification to extract the desired
18 eye-frames within the input image. We then use multi-stage convolution to find the optimal
19 horizontal and vertical coordinates of the pupil within the identified eye-frames. For this purpose,
20 we define an adaptive kernel to deal with the varying resolution and size of input images.
21 Furthermore, a dynamic threshold is calculated for reliable identification of the best-matched
22 candidate. We evaluated our method using various statistical and standard metrics along-with a
23 standardized distance metric we introduce first time in this study. Proposed method outperforms
24 previous works in terms of accuracy and reliability when benchmarked on multiple standard
25 datasets. The work has diverse artificial intelligence and industrial applications including human
26 computer interfaces, emotion recognition, psychological profiling, healthcare and automated
27 deception detection.

28 **Keywords:** Pupil detection; Deep eye, Iris detection; Eye centre localisation; Eye gaze; Facial
29 analysis, Image convolution; Machine intelligence, Pupil segmentation
30

## 1. Introduction

32     Detection and localization of the objects within images or real time video frames is considered
33 an essential task in various computer vision algorithms [1]. Various studies have addressed the
34 detection and tracking of facial landmarks including the iris and pupil which has various
35 applications particularly, eye gaze estimation for human-machine interfaces. Control of assistive
36 devices for disability [2], driver safety improvements [3-4], the design of diagnostic tools for brain
37 diseases [5], cognitive research [6], automated deception detection system (ADDS) [7] and academic
38 performance analysis [8] are some examples of such applications.
39     Research studies for the eye detection and eye tracking mostly focus on the iris and pupil
40 localization. Once the coordinates of pupils are determined, it can be used for the eye tracking, gaze
41 estimation and eye movements within the images and video frames [6]. Eye images can be
42 characterized by the intensity distribution of iris, pupil and the cornea, in addition to their shapes. It
43 should be noted that various aspects can influence the appearance of the eye including the viewing

44  angle, ethnicity, head position, eye colour, light conditions as well as the texture, eye state (e.g. half
45  closed, fully closed) and current wellbeing [6].
46      Overall, eye detection techniques can be classified as shape-based, feature-based,
47  appearance-based and hybrid methods. In the shape-based methods, open eyes are described by
48  their shapes including the pupil and iris contours as well as shape of the eyelids [9-11]. For the
49  feature-based methods, objective is to identify the local features within the eye that are less sensitive
50  to the varying illumination as well as viewpoint [12-15]. Appearance-based methods depend upon
51  detecting and tracking of the eyes using the photometric look which is characterized by colour
52  distribution and filter responses to eyes and their surroundings [16-18]. The hybrid methods aiming
53  to combine various techniques to mitigate the particular disadvantages of these methods [19-20].
54      Standard methods in gaze estimation are based on corneal reflections that needs an accurate
55  localization of the pupil centre as well as the glints [21]. Pupil and glints localization algorithms are
56  usually based on image processing such as morphological operators for the detection of contour [22]
57  and intensity threshold identification followed by the fitting using ray-based ellipse [23].
58  Topography based hybrid method is introduced in [24] which uses series of filters for the iris centre
59  estimation. However, these techniques assume that the pupil exists in the darkest area of the input
60  image and may susceptible to varying illumination conditions that might require manual tweaking
61  to the threshold parameters [25].

62                              Table 1. Eye movements and classification algorithms

| Reference | Model | Aims and Feature Used |
|---|---|---|
| [26] | Hidden Markov model | Use of fixation count, fixation durations to distinguish between expert and novice participants |
| [27] | Multi-layer perceptron (MLP) | Use pupil size & point-of-gaze for predicting the users' behaviours (e.g., word searching, question answering, looking for the most interesting title in a list) |
| [28] | Naïve Bayes classifier | Use of fixation duration, mean and standard deviation to identify various visual activities (e.g., reading, scene search) |
| [29] | MLP | Use of Pupil dilation, gaze dispersion to classify various tasks on decision making |
| [30] | Decision tree, MLP , support vector machines (SVM), linear regression | Use of fixation rate, fixation duration, fixations per trial, saccade amplitude, relative saccade angles to identify eye movements to predict visualization tasks |

63      There are four main eye movement behaviours which are likely to show different details related
64  to cognitive efforts when responding to tasks including blinks, pupillary responses, fixations, and
65  saccades [31]. Blinking represents the involuntary deed of opening and closing the eyelids. Pupillary
66  responses are the changes in pupil size restrained by the involuntary nervous system. Fixation
67  represents the collection of gaze points that are relatively stable and near in spatial and temporal
68  vicinity. Saccade represents the rapid and small eye movements when moving from one object to
69  another [31]. These four eye-movement behaviours reveal the details about cognitive efforts and
70  therefore can be used as suitable inputs for designing the machine learning (ML) systems as
71  illustrated in Table 1 which shows various supervised ML algorithms to predict categorical
72  responses from the eye movements.
73      In addition to conventional methods, existing works also utilize the deep learning (DL)
74  approaches for the pupil detection while using hierarchical image patterns to enhance and eliminate
75  artefacts with Convolutional Neural Networks (CNNs). For instance, [21] proposed the use of fully
76  connected CNNs for segmentation of the entire pupil area in which they trained the network on 3946
77  video oscillography images. These images were hand annotated and generated within a laboratory
78  environment. The authors claim that the proposed network enables them to perform elliptical
79  contour detection, pupil centre estimation and blink detection. More explicitly, pupil centre are
80  predicted with a median accuracy of one pixel and gaze estimation accuracy within 0.5 degrees.
81  However, varying image resolution might provide different accuracy measures. More specifically,

82  [32] indicated the eye tracking as an important tool that can have a range of applications from
83  scientific research to commercial sector. The authors show that the use of tracking software based on
84  commodity hardware including tablets and smartphones, allows these advanced technologies to be
85  available for everyone. The system is called iTracker which uses a CNNs model indicating 2.53cm
86  and 1.71cm prediction error without calibration on tablets and smartphones respectively which is
87  reduced to 2.12cm and 1.34cm using calibration.
88  Research presented in [23] proposed a pipeline of two CNNs cascaded for pupil detection.
89  Authors claim that their method outperforms state-of-the-art techniques with detection rate up to
90  25% while avoiding computational complexity. To benchmark their proposed technique, 79000 hand
91  labelled images were used in which 41000 were complementary to existing images from the
92  literature. A similar work is presented in Naqvi et al. [33] which indicate that automobile accident
93  deaths could be minimized using drivers' gaze region to provide their point of attentions. In this
94  respect, the authors suggest the use of DL for gaze detection with the use of near-infrared camera
95  sensors. They incorporate driver head and eye movement into their study. Gaze estimation accuracy
96  was benchmarked using loosely correct estimation rate and strictly correct estimation rate in which
97  the study claim achieving good accuracy when benchmarked with the previous gaze classification
98  techniques.
99  Recent work that uses the CNNs based deep learning model for the pupil estimation [34]
100 indicate around 70% accurate estimations while error threshold is within the 5 pixels. However, this
101 accuracy is limited to be used in real time specifically, the applications that consider
102 micro-movements within the eyes such as ADDS [7]. Similar work that uses CNNs for the pupil
103 detection [35] indicates varying detection rate (70-90%) with respect to the tolerance level as pixel
104 error and dataset they employed for testing. The study outcomes clearly indicate the trade-off
105 between the error tolerance level and accuracy measure. Furthermore, the performance metric used
106 in these studies is not standard (i.e. the error as number of pixels) and might produces varying
107 accuracy with respect to image size and resolution. In contrast to CNNs, [36] utilizes the wavelet
108 transform to extract the distinguishing features while SVM is used for the pupil classification. This
109 work indicates 88.79% of accurate pupil estimation on a benchmarked dataset while utilizing the
110 standard validation metric.
111 Despite the variety of existing methods for the pupil localisation, further improvements are
112 required in terms of a precise estimation for the pupil location. For instance, DL-based pupil
113 localisation and gaze estimation in [21] uses pixel distance to validate the performance which is not a
114 standard representation of the error in case of varying resolutions. Furthermore, the validation is
115 performed on a dataset containing artificially rendered images which in most cases, does not reflect
116 the real time dynamics. Likewise, [37] presented gaze estimation that utilizes the DL-based facial
117 landmarks detection following the image segmentation to identify the pupil within the input
118 images. However, the 81% accuracy produced by the algorithm on a benchmark dataset indicates
119 the lack of preciseness in pupil localisation that might lead to the incorrect gaze estimation.
120 Furthermore, this study along with [23, 24] utilizes a static threshold while considering the pupil as
121 the darkest area within the image that may susceptible to various illumination conditions [25] and
122 low-resolution images. Likewise, the use of static size kernel for the template matching to find out
123 the best-matched candidate (i.e. pupil in this case) within the image might causes local maxima. For
124 instance, a smaller sized kernel may cause attention to noisy details (i.e. local maxima) whereas,
125 larger size may lead to mismatches and incorrect estimation of pupil location [38].
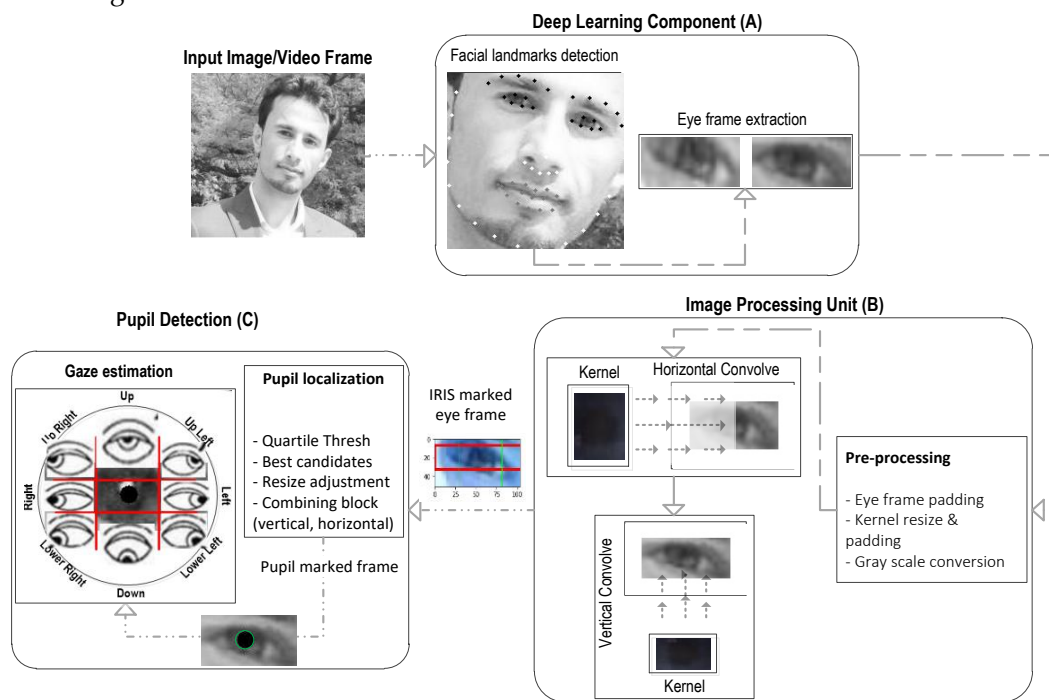126 In proposed work, we introduce an efficient algorithm for the pupil identification within
127 low-resolution images (and video frames) using a composite of DL and image processing
128 algorithms. To clarify the novelty of this paper, the contributions are outlined as follows. a) utilizing
129 the pre-trained DL model to identify the facial landmarks and extraction of desired eye-frames
130 within the input images; b) unidirectional cascades of two-dimensional (2D) convolution is used to
131 determine the pupil coordinates within the eye-frames of varying characteristics; c) an adaptive size
132 of kernel is used to deal with the varying size of input images (i.e. eye-frame) during the template
133 matching; d) we used a dynamic threshold to identify the best matched candidate more reliably; e)

134 for the first time, we introduce a relative error metric to measure the standardized distance (i.e.
135 error) between the estimated and actual pupil centres; f) we validated the proposed methodology
136 over multiple publicly available and benchmark datasets containing high diversity in gaze positions,
137 participants background, lighting illuminations, image background, and comparatively smaller size
138 of eye-frames.
139     The remainder of this paper is organized as follows. Section 2 entails the proposed
140 methodology and algorithms. Section 3 presents the detailed experimental design and newly
141 introduced evaluation metric. Statistical results and technical discussions are presented in Section 4
142 followed by a conclusion and future works in Section 5.

## 2. Proposed Method

144     The proposed pupil detection utilizes composite of techniques along with new algorithms while
145 leveraging the DL-based facial landmark detection [39] to extract the eye information within an
146 image/video frame. Existence of background noise and dark patches within the image frame and
147 specifically prominent eyebrow parts, are normally detected as pits that might cause mismatch for
148 computer vision-based iris and pupil detection [24, 38, 40]. However, this issue can be resolved
149 readily by utilizing modern DL algorithms for a reliable face and eye-frame extraction from an
150 ordinary quality images or video frames. In the first step, we utilize the facial landmark detection to
151 extract the desired segments containing only the eye-frames (both left and right) from input image.
152 We then convolve the extracted eye-frame with a pre-defined kernel in horizontal and vertical
153 directions to identify the iris and pupil respectively within the eye-frame. We adapt the kernel size
154 dynamically with respect to the varying eye-frame size to resolve the possible occurrences of local
155 maxima being false representation of best matched patches. We further define a dynamic threshold
156 for the identification of best-matched patch within the current eye-frame to reduce the impact of
157 noisy matches. Figure 1 shows the sequential processing in our work to identify the pupil
158 coordinates within an input image/video frame. The major components are: a) DL-based eye-frame
159 extraction, b) image processing based iris localisation, and c) pupil detection, which are detailed in
160 the following sub-sections.



161

**Figure 1.** Sequential processing components of the proposed method comprising A) DL library (i.e. Dlib-ml) for the eye-frame extraction, B) computer vision algorithm for localizing the potential iris and pupil candidates within eye-frames, C) post-processing for the pupil coordinate measurement. In images, eyes view is reversed (e.g. the left eye in an image is the right eye in actual and vice versa)

166 *2.1. Eye frame extraction*

167     The DL component utilises a well-known toolkit (Dlib-ml) [39] which can reliably identify the
168 facial landmarks while producing extensive fiducial points (68 in total) on the face including eye
169 corners and eye lids as shown in figure 1(A). We first extract the face rectangle from an image using
170 Dlib-ml that not only removes the unnecessary portion of input frame but also helps to eliminate the
171 major noisy components that might exist in background region of the image frame. Within the face
172 region, we then note the identified extreme points (left, right, top, bottom) for eye corners and eye
173 lids which are used to crop the exact eye-frames within the identified face rectangle. This is one of
174 the major advantages of using Dlib-ml which reliably eliminates the unnecessary portion of an
175 image and extract the exact region of interest (i.e. eye-frames in this case) from the input frame. Only
176 the input images (or video frames) with exactly one face rectangle and two eye-frames are
177 considered as 'valid'. The output of this component in form of eye-frames (left, right) are processed
178 further to identify the iris and pupil within the image.

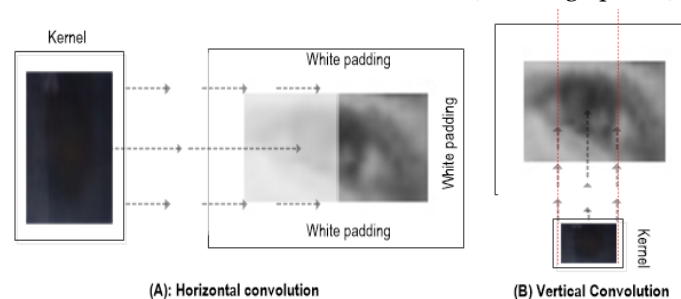179 *2.2. Iris segmentation and pupil localization*

180     Following the eye-frame extraction, a convolution function is applied for the template matching
181 between a custom kernel and eye-frame to localise the best matching segment within the eye-frame.
182 Firstly, we built a custom kernel representing 100 iris frames (cropped from eyes frames) randomly
183 chosen from datasets described in Section 3. The advantage of custom kernel over an ordinary black
184 colour kernel, is a more generalized representation of an iris for a diverse population and
185 morphology characteristics (e.g. geometry, patterns within the iris, colour etc.). Another common
186 factor that can affect the template matching performance, is size of the template (i.e. kernel). Smaller
187 sized kernel may cause attention to noisy details (i.e. local maxima) whereas, larger size may lead
188 mismatches and incorrect estimation [38].
189     To resolve this issue, the adaptive size kernel is employed using the interpolation and
190 extrapolation techniques where the size ($w_k \times h_k$) varies with respect to the input frame size (i.e.
191 eye-frame). Furthermore, eye-frame (*E*) is padded with a rim of white pixels (see figure 1 and figure
192 2) to enlarge it enough that the convolution kernel (*K*) fits inside the padded image to provide all
193 possible best matches identification (i.e. between kernel *K* and overlapped eye-frame patches of a
194 size similar to *K*). More specifically, when the desired patch (i.e. iris) is located at extreme positions
195 (e.g. looking extreme left/right positions).

$$y[i,j] = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} K[m,n].\,E[i-m, j-n] \qquad (1)$$

196

197     Equation 1 represents a 2D convolution function where *E* is the current eye-frame (within the
198 input image) to be convolved with the kernel matrix *K* resulting *y* as the output image. The indices *i*,
199 *j* and *m, n* represent the indices within the *E* and *K* matrices (i.e. image pixels), respectively.



(A): Horizontal convolution      (B) Vertical Convolution

200

201     **Figure 2.** Horizontal convolution (A) and vertical convolution (B) between adaptive size kernel *K*
202     and white outlined eye-frame *E*.

203     In contrast to the ordinary way of 2D-convolution where kernel *K* slides along *E* with a fixed
204 overlapping window (usually 1 pixel) in both horizontal and vertical directions, we perform a
205 comparatively simple and efficient convolutional steps (only one slide per horizontal and vertical

206    directions) as shown in figure 2. The reason behind an adaptive kernel selection is the geometric
207    features of iris and pupil which are considered approximately circular and black compared to the
208    rest of the eye with pupil as the most dark segment. First, kernel height $h_k$ is resized to eye-frame
209    height (i.e. $h_e = h_k$) and width $w_k$ is set to 0.4 of the eye-frame width. The convolution function then
210    slides through $E$ in the horizontal direction to determine the x-coordinate of iris centre within the $E$.
211    It compares the overlapped patches of $E$ ($w_k \times h_k$) against $K$ to calculate the matching scores at each
212    horizontal stride (i.e. 1 pixel). The normalised correlation coefficient calculates a total matching score
213    for the current patch in $E$ using equation 2.

$$S(x, y) = \frac{\sum_{x'y'}(K'(x', y') \cdot E'(x + x', y + y'))}{\sqrt{\sum_{x'y'} K'(x', y')^2 \cdot \sum_{x'y'} E'(x + x', y + y')^2}} \quad (2)$$

214

215    Where $S(x, y)$ is the matching score of current overlap $(x, y)$ between $K$ and $E$ patch of size
216    equal to $K$ ($w_k \times h_k$). The summation in equation 2 is performed over the $K$ and $E$ patch where $x' =$
217    $0...w_{k-1}$, $y' = 0... h_{k-1}$. As the kernel height $h_k$ is aligned with height of the eye-frame (i.e. $h_e = h_k$), there are
218    no vertical overlapping (i.e. no vertical overlapping/strides) which means, the kernel will only be
219    able to move along $E$ in the horizontal direction while computing the matching scores for
220    overlapped patches in $E$.
221    Once all the horizontal matching scores are calculated, the next step is to find the coordinates of
222    the best matching segment. There have been several approaches to select the optimal match but the
223    candidate with maximum match have been commonly used in similar works [12, 38, 41]. However, it
224    can easily cause local maxima specifically in low-resolution images [38]. Likewise, using a
225    predefined matching threshold can provide varying matching scores regarding the environment and
226    can also mislead because of varying dynamics such as illuminations. We utilised quantile measure to
227    select all candidates ($M$) in the horizontal direction) that crosses the adaptive threshold of 90th
228    percentile    of    the    matching    scores    sorted    in    ascending    order.    i.e.
229    $M \in SC_h$ such that $\forall SC_h > 90^{th}$ percentile of sorted $SC_h$.

230    The mean of horizontal (*x-axis*) coordinate of $M$ selected patches is calculated using (3) which
231    represents the x-coordinate of top-left corner ($R_{xy}$) of the final best matched patch (i.e. estimated iris
232    rectangle).

$$R_x = \frac{\sum_{i=1}^{m} M_x(i)}{m} \quad (3)$$

234    Where $m$, are the total number of elements (i.e. best-matched candidates) in $M$, $M_x$ is the
235    horizontal coordinate of corresponding best-matched candidates $M$.

236    The iris rectangle $I$ is identified using $R_x$ and kernel width $w_k$ which is then used for the vertical
237    convolution to identify the y-coordinate of iris centre. Similar to horizontal convolution-based
238    matching, kernel height $w_h$ is resized to 0.4 of the height of $I$ for overlapped stride matchings while
239    keeping the width same. Vertical convolution steps are then performed to compute the matching
240    score for $K$ and overlapped patches of $I$ along the vertical direction only. The output matrices $SC_v$
241    contains all the corresponding matching scores for vertical convolutions between the $K$ and $I$
242    overlapped patches. The quantile measure is used in a similar way to select all candidates ($N$) in the
243    vertical direction) that crosses the adaptive threshold of 90th percentile of the matching scores sorted
244    in ascending order where; $N \in SC_v$ such that $\forall SC_v > 90^{th}$ percentile of sorted $SC_v$. The mean of
245    vertical (*y-axis*) coordinate of $N$ selected patches is then calculated using equation 4 which represents
246    the y-coordinate of top-left corner ($R_{xy}$) of the final best-matched patch (i.e. estimated pupil
247    rectangle).

$$R_y = \frac{\sum_{i=1}^{n} N_y(i)}{n} \quad (4)$$

248

249  where *n*, is the total number of elements (i.e. best-matched candidates) in *N*, $N_y$ is the vertical
250  coordinate of corresponding best-matched candidates *N*.

251  Finally, the centre coordinates of the best-matched patches within *E* in horizontal ($C_x$) and
252  vertical directions ($C_y$) represent the pupil location along the *x-axis* and *y-axis* respectively and are
253  calculated as:

$$C_x = R_x + w_k/2, \ C_y = R_y + h_k/2 \tag{5}$$

254

255  where $w_k$, $h_k$ are the width and height of kernel *K*, respectively. Algorithm 1 summarizes all the
256  sequential steps involved in the proposed methodology to determine the pupil coordinate within an
257  image frame.

258

**Algorithm 1:** Proposed algorithm for iris detection and pupil localization in an image/video frame

**Inputs**: image/video frame *F*, a custom-defined kernel frame *K*
**Output**: Pupil coordinates (*Cx, Cy*), iris rectangle (top-left; bottom-right)

**STEP1:**
- Initialise validation *Score = 0* for current *F*
- Use *Dlib-ml* for the facial landmark detection within input frame *F*
- Crop the face rectangle (*Face*) using the detected landmarks
- IF count (*Face*) ==1 (i.e. exactly one face in image is found)
    - *Score* ++
    - Extract the eye-Frames ($E_L$, $E_R$) for *left* and *right* eye
    - IF count ($E_L$, $E_R$) ==2. i.e. exactly 2 eyes within *Face* rectangle
        - *Score* ++
        - *Goto* STEP 2
    - ELSE
        - Mark it as invalid frame
        - *Goto* STEP 1 for the next *F*
- ELSE
    - Mark it as invalid frame
    - *Goto* STEP 1 for the next *F*

**STEP2:**
- Foreach *eye-frame E* in $E_L$, $E_R$
    - Convert *E* into grayscale
    - Outline *E* with white paddings
    - Adapt the kernel *K height* to *height* of *E* and *width* to 0.4**width(E)*
    - Convolve *K* with *E* by sliding *Horizontally* with 1-pixel stride/sliding window
    - Store the matching scores for overlapped *E* patches in a vector $SC_h$
    - Store the horizontal elements with high matching scores in lists *M* for

        $M \in SC_h$ such that $\forall\, SC_h > 90^{th}$ percentile of sorted $SC_h$.

    - Find the *top-left* of best-identified iris rectangle by taking mean ($\mu$) of x-coordinates

        for *M* (i.e. $R_x$) using equation 3

    - Find the iris rectangle *I*, using $R_x$ and $w_k$
    - *Goto* STEP3
- End Loop

**STEP3:**
- Adapt the kernel *K height* to 0.4**height(I)* for vertical convolution
- Convolve *K* with *I* by sliding *Vertically* with 1-pixel stride/window
- Store the matching scores for overlapped *I* patches, in a vector $SC_v$
- Find the elements with high matching scores (call them *N*) where

    $N \in SC_v$ such that $\forall\, SC_v > 90^{th}$ percentile of sorted $SC_v$.

- Find the *top-left* coordinate of best-identified rectangle by taking mean ($\mu$) of y-coordinates of

    *N* (i.e. $R_y$) using equation 4

- Find the pupil centre *Cx, Cy* by adding width and height of *K* into $R_x$ and $R_y$ respectively
    using equation 5.

### 3. Experimental Design

We conducted detailed experiments to validate the proposed methodology while using various datasets and validation metrics. We also performed a critical analysis based on various conditions and validated the proposed algorithm while considering the diversity in validation datasets as well as validation metrics. Following sections explain the validation datasets and metrics along-with detailed experimental design.

*3.1. Datasets*

To validate the proposed methodology and reliable performance measure, we used three different publicly available datasets. The first dataset is known as Talking-Face [42] and have been used in previous works [37]. This dataset contains 5000 video frames captured during the engaged conversation from a person for 200 seconds. The original objective of this dataset was to model the facial behaviour during a natural conversation. Data is captured with a static positioned camera with a frame size of 720x576 pixel. Every frame is annotated semi-automated manner containing 68 facial points including the pupil coordinates. Following our validation check in Algorithm I (i.e. frames with exactly 2 eyes/frame) and removing the fully closed eyes (manually, found 280 images) images, we are left with 4720 frames for the validation purpose. The dataset contains varying gaze positions, facial and body movements, diverse natural expressions and variations in eye-state (e.g. closed, open, half closed) . However, because it is captured from individual person , the diversity within the eye characteristics is very limited. In other words, there are no variations in terms of eye characteristics (e.g. Iris or pupil colour, intensity, iris pattern etc.) and hence, not very challenging for the algorithm validation.

In contrast to Talking-Face, we used the BIO-ID dataset [43] which is comparatively more challenging and has been used as a benchmark in various relevant studies such as [37, 41]. The data was acquired from 23 different subjects during multiple sessions and has 1521 images in total containing varying gaze positions, illuminations, background scene, eye features (e.g. eye colour, gender, ethnicity, iris size), camera focus and hence eye-frame (and face rectangle) size. The interesting aspect of this dataset is a comparatively lower resolution (grayscale 384x288 pixel) that makes the validation of pupil localisation algorithm more challenging but reliable. Besides, the dataset contains natural expressions such as images with half-closed eyes that further help to measure the validity of the proposed algorithm. Our algorithm detects only seven frames as invalid (i.e. not containing exactly two eyes) whereas we found 45 images (manually) with fully closed eyes that were excluded, resulting 1469 remaining dataset for validation purpose.

Furthermore, we evaluated our method on comparatively larger dataset known as GI4E [44] containing more diversity involving various morphology types (e.g. eye size, eye/iris features, gender, ethnicity, varying background and illuminations). It should be noted that despite higher resolution images (800×600 pixels), size of the eye-frame rectangles is comparatively small. This is because of the larger distance of the capturing device from the subject resulting lower ratio of eye-frame to entire image. In other words, the whole frame covers more background pixels as compared to the actual face within the image which makes the eye-frame and hence iris/pupil localization more challenging. The dataset is much diverse containing 103 subjects (each with 12 images) with 1236 total images involving 12 different gaze position. Also, there is no open eyes or invalid frame in this dataset.

*3.2. Validation Metrics*

One of the important factors in validation of the pupil detection and proposed work is the metric we chose for the performance measure. This is because of the nature of pupil localization problem. For instance, the absolute error in the estimated pupil/eye centre and actual eye centre might vary with respect to image size/resolution. Hence the standard distance measure such as Euclidean distance (ED) and/or $R^2$ coefficient will not give a true representation of the accuracy measure. The authors in [43] introduced a relative error measure ($d_{eye}$) to deal with this issue which

308  has been utilized in various related works [37, 38, 43, 45]. It uses the maximum of the estimated pupil
309  coordinates distances from left and right eyes $(d_l)$ and $(d_r)$ respectively, between the actual eye
310  centres $(C_l,\ C_r)$ and the estimated ones $(\tilde{C}_l,\ \tilde{C}_r)$ using equation 6.

311
$$d_{eye} = wec = \frac{max\ (\ \|\tilde{C}_l - C_l\|,\ \ \|\tilde{C}_r - C_r\|\ )}{\|C_l - C_r\|} \tag{6}$$

312  For the normalisation, the calculated distance is divided by the distance between two actual eye
313  centres $\|C_l - C_r\|$ as shown in equation 6. The normalisation factor makes the error measure
314  independent of the image scale and hence eye-frame size. Furthermore, [37] used best eye centre
315  (*bec*) which utilizes the minimum of the error between estimated and actual centres as:

316
$$d_{eye} = bec = \frac{min\ (\ \|\tilde{C}_l - C_l\|,\ \ \|\tilde{C}_r - C_r\|\ )}{\|C_l - C_r\|} \tag{7}$$

317  Although the *wec* (i.e. worst eye centre) metric provides a relative error estimate, it is based on
318  some assumptions such as '*on average population, the distance between the inner eye corners is equal to*
319  *width of a single eye of the corresponding subject*'. Likewise, a relative error of $d_{eye}$= 0.25 is considered as
320  half of an eye width which may not be valid in every case. Interested readers can get further details
321  in [43] study.
322  To further deal with the metric generalisation issue, we first time introduce a standardized
323  error measure (*S_{ED}*) as a function of distance between the estimated and actual coordinates within an
324  eye-frame. It calculates the relative distance as percentage of the total possible ED (i.e. error)
325  between the actual and estimated pupil coordinates. The *S_{ED}* measure interprets the error within the
326  single eye-frame without depending on the second eye or interpupillary distance used in other
327  related works. Besides, the *S_{ED}* metric can measure the relative error regardless of image/face or
328  eye-frame size and hence the image resolution. Mathematically, the proposed *S_{ED}* is defined as:

329
$$S_{ED} = \frac{\sqrt{(Cx_e - Cx_a)^2 + (Cy_e - Cy_a)^2}}{\sqrt{(x_{min} - x_{max})^2 + (y_{min} - y_{max})^2}} \times 100 \tag{8}$$

330  Where $Cx_e, Cx_a$ represent the estimated and actual pupil horizontal coordinates respectively
331  and $Cy_e, Cy_a$ represent the estimated and actual pupil vertical coordinates respectively. The
332  $x_{min}, y_{min}$ are coordinates of the nearest corner of eye-frame (usually top left corner) whereas,
333  $x_{max}, y_{max}$ are coordinates of the farthest corner of eye-frame (usually bottom right). The numerator
334  in equation 8 represents the error (in terms of pixels) between the actual and estimated positions
335  whereas the denominator is the total possible error and is used as a normalisation factor. The
336  resulting *S_{ED}* gives the percentage error representing a standardised distance between actual and
337  estimated pupil positions in pixels which is not affected by the image size and resolution. In addition
338  to evaluate the pupil detection techniques, the proposed standardised distance measure can also be
339  useful for other related works such as object localisation, image segmentation and object tracking
340  etc.
341  In summary, a comprehensive comparative analysis is performed to evaluate the proposed
342  methodology using aforementioned metrics including *wec*, *bec*, and *S_{ED}* along with other standard
343  accuracy measures including the ED, absolute mean difference, and $R^2$ (coefficient of determination).

344  **4.  Results and Discussions**

345  Following the experimental design, performance of the proposed pupil detection approach is
346  evaluated using various gold standards, validation metrics and benchmarked datasets. As discussed
347  in the experimental design, it is important to use appropriate evaluation methods due to nature of
348  the problem. To maintain the reliability in our performance measure, we utilised different metrics as
349  well as the newly introduced *S_{ED}* in this work.
350

351    **Table 2.** Performance analysis of the proposed model using *wec*, *bec* with varying error threshold

| Dataset | *wec*(%) | | *bec*(%) | |
|---|---|---|---|---|
| | Error ≤ 0.05 | Error ≤ 0.1 | Error ≤ 0.05 | Error ≤ 0.1 |
| BIO-ID | 94.5 | 100 | 98.34 | 100 |
| Talking face | 97.10 | 100 | 99.7 | 100 |
| GI4E | 95.05 | 100 | 98.71 | 100 |

352    Table 2 summarizes the results achieved from the proposed approach using *wec* and *bec* metrics
353    that have been used in recent similar works [36-38, 43-45]. We are specifically interested in *wec*
354    measure when *error≤0.05* which indicates the model estimation within the *pupil* diameter (i.e. more
355    restricted). Best accuracy achieved by the proposed method is 97.1% while tested over the
356    *Talking-Face* dataset which outperforms the 89.59% presented in recent work[37] that uses the same
357    dataset. The high accuracy is expected because of the comparatively less challenging nature of
358    dataset (see Section 3.1). Firstly, the dataset contains high resolution images. Secondly, the data is
359    captured from only one person hence, a generalization of *iris* and eye pattern is easily detected. It is
360    important to note that despite the dataset is collected from single person, it contains high variations
361    in terms of gaze, head movements, facial expressions and sufficient quantity (i.e. 5000 images) with
362    annotated pupil coordinates. On the other hand, the proposed method achieves 100% *wec* accuracy
363    while tested for error threshold≤0.1 indicating the robustness of the proposed methodology. This
364    means that model estimation about *pupil* coordinates are within the *iris* in all cases (i.e. 5000 images).
365    Overall, proposed method outperforms the most recent works related to pupil localization [37]
366    while evaluated on the *Talking-Face* dataset.



368    **Figure 3.** Comparison of estimated pupil coordinates using proposed model, with actual annotated
369    coordinates (BIO-ID dataset) using R-squared error

370    To further evaluate the model performance, the BIO-ID dataset is used which contains various
371    subjects, high variations in gaze, head pose and body movements. Furthermore, the image quality
372    (i.e. resolution) is comparatively lower (i.e. 286x384) which makes it more challenging when
373    focusing the identified eye-frame and/or iris/pupil within the image. Also, a large proportion of the
374    entire image contains background rather than the face itself which makes the dataset further
375    challenging as addressed by the previous works [38]. Despite the associated challenges, proposed
376    approach shows robust pupil estimations as shown in Table 2. The model indicated significant
377    improvements with 94.5% *wec* measure with error threshold≤0.05 when benchmarked with the
378    works of [37] and [40] of 81% and 84%, respectively. Furthermore, the model indicated 100%
379    accuracy when evaluated for error threshold≤0.1 which means that pupil localization is within the

380  iris in all cases (i.e. 1521 cases in total). Despite the 100% of *wec and bec* accuracy for error
381  threshold≤0.1, the main focus is to maximize the *wec* accuracy (which is the most challenging) with
382  minimum error threshold (i.e. ≤0.05) to restrict the model estimation within the pupil diameter.
383      Figure 3 shows the $R^2$ coefficient for the proposed model tested on BIO-ID dataset. It can be
384  observed that x-axis and y-axis estimated coordinates are almost overlapping to actual annotations
385  with $R^2$ value of 0.993 and 0.998 for x-axis and y-axis respectively. Although, $R^2$ is a well-known
386  statistical measure to determine the goodness of model fit, it might not be effective for validating the
387  model estimation in pupil detection or similar problems because of the varying error rate with
388  respect to the image size (and resolution).

389      **Table 3**. Performance comparison between previous works based on *wec* measure using BIO-ID
390                                                          dataset

| *wec % accuracy with varying error (e) threshold* | | | |
|---|---|---|---|
| **Methods** | *e<0.05* | *e<0.1* | *e<0.15* | *e<0.2* |
| [24] | 81.1 | 94.2 | 96.5 | 98.5 |
| [36] | 88.7 | 95.2 | 96.9 | 97.8 |
| [37] | 80.9 | 91.4 | 93.5 | 96.1 |
| [38] | 82.5 | 93.4 | 95.2 | 96.4 |
| [40] | 84.1 | 90.9 | 93.8 | 97.0 |
| [41] | 57.2 | 96.0 | 98.1 | 98.2 |
| [43] | 38.0 | 78.8 | 84.7 | 87.2 |
| [45] | 47.0 | 86.0 | 89.0 | 93.0 |
| [46] | 85.8 | 94.3 | 96.6 | 98.1 |
| **Proposed Model** | **94.5** | **100** | **100** | **100** |

391      Table 3 summarizes the comparative results from various previous works while weighted over
392  the challenging BIO-ID dataset using *wec* metric with varying thresholds. It can be noticed that the
393  proposed model outperforms (94.5%) all previous works specifically with the most restricted error
394  threshold≤0.05. Recent works that uses similar approach [37] achieved an accuracy of 80.9% and
395  82.5% [38] with e≤0.05 whereas best accuracy of 88.79% is indicated by [36] that are significantly
396  lower than the proposed method. Research study in [21] presented a robust technique for the pupil
397  localization and gaze estimation, however, the measured performance is not standard (i.e. uses the
398  mode of pixel distance which is not the true representation of error with varying resolutions).
399  Furthermore, the validation is performed on different dataset containing artificially rendered images
400  which in most cases, does not reflect the real time dynamics.
401      Besides the Talking-Face and BIO-ID datasets, we evaluated the performance of proposed
402  approach on another challenging dataset GI4E. It can be noted from Table 2 that our model produces
403  95.05% *wec* and 98.71% *bec* accuracy respectively with critical threshold≤0.05. While most of the
404  existing works used BIO-ID as benchmark dataset, some of them also used GI4E to evaluate their
405  techniques. For instance, recently study on eye centre localisation [24] reported 93.9% *wec* accuracy
406  on GI4E dataset which is slightly lower than our approach (i.e. 95.05%) however, their accuracy was
407  decreased to 881.2% when tested on BIO-ID dataset. This indicates the robustness of proposed
408  approach for pupil detection in varying datasets containing diversity in terms of eye colour, gaze
409  position, facial emotions and real movements. Similarly, [46] indicated 89.28% *wec* on GI4E dataset
410  which are significantly lower than the proposed approach. A clustering-based approach [47]
411  produced mean pixel error of 2.73 pixels as compared to proposed model with 1.7 pixels while
412  validated on GI4E. However, it is important to be noted that this metric does not represent a
413  standard accuracy measure as described in Section 3.2.
414      In addition to *wec*, [24, 41] used average point-to-point error ($m_{e17}$) with the inter-ocular distance
415  between the left and right eye pupil. Recent works [21] that uses the DL to localize the pupil and
416  estimate the gaze position also employed the median of absolute difference in x-axis and y-axis.
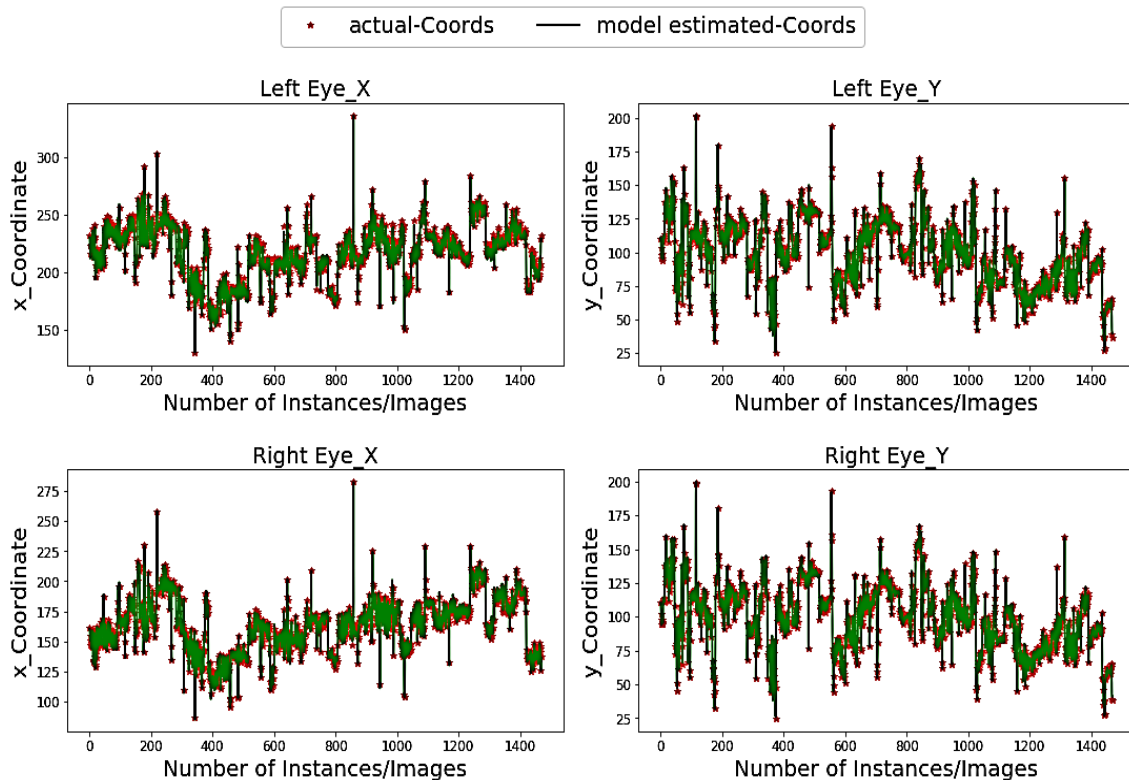
417  However, variations in image size, zoom-in/out due to body/head movements and/or camera
418  positions might affect the mean difference in corresponding error estimate resulting variations in
419  accuracy measure. The *wec* metric which has been used extensively in related works such as [24,
420  36-38, 40, 45, 46], gives a comparatively better indication of the performance measure. However,
421  these metrics measures the performance in terms of coordinate estimation within the pupil/iris
422  diameter with a varying error threshold as shown in Table 2. Also, it is based on relative error
423  assumption ($d_{eye}$= 0.25) as half an eye width, which may not be true in every case. Therefore, model
424  estimations and performance (specifically in pupil localization task) is needed to be evaluated using
425  more standard metric representing the distance between estimated and actual pupil coordinates.

426  **Table 4.** Comparing model estimations using newly introduced S$_{ED}$, Euclidean distance (ED), R$^2$, and
427  absolute error metrics

| Dataset | $\mu|x_a-x_e|$ | $\mu|y_a-y_e|$ | R$^2$_x | R$^2$_y | ED($c_a$, $c_e$) | %ED($c_a$, $c_e$) |
|---------|------|------|------|------|------|------|
| BIO-ID | 1.04 | 0.57 | 0.993 | 0.998 | 1.43 | 3.98 |
| Talking face | 1.23 | 0.97 | 0.990 | 0.956 | 1.96 | 2.49 |
| GI4E | 1.32 | 0.71 | 0.996 | 0.999 | 1.70 | 3.87 |

428  To overcome this issue, we first time introduce a standardized Euclidean distance (S$_{ED}$) which
429  represents the percentage distance error as ED using equation 7 (see Section 3.2). The error
430  represents the displacement between the actual and estimated pupil coordinates as a percentage of
431  the whole image size (i.e. eye-frame) in terms of number of pixels. The major advantage of S$_{ED}$ is a
432  standard representation of the error which can be used to measure the accuracy regardless of image
433  size and resolution which is not the case in *wec*, m$_{e17}$ and other metrics used in most of the existing
434  studies. Table 4 presents the comparative analysis of proposed model estimations in terms of mean
435  pixel difference in each axis, for both eyes (left and right), R$^2$ coefficient, ED between centre of
436  estimated and actual pupil coordinates and the newly introduced S$_{ED}$. The proposed method
437  indicates 1.04 and 0.57 absolute pixel error on x-axis and y-axis respectively (i.e. 0.8 on average for
438  both) as compared to 2.91 in [47] on BIO-ID dataset. Similarly, a DL-based model in [35] indicated
439  their optimal performance with pixel error>10. However, they used different datasets which in case
440  of high resolution, is not comparable with proposed method and clearly indicates the need of
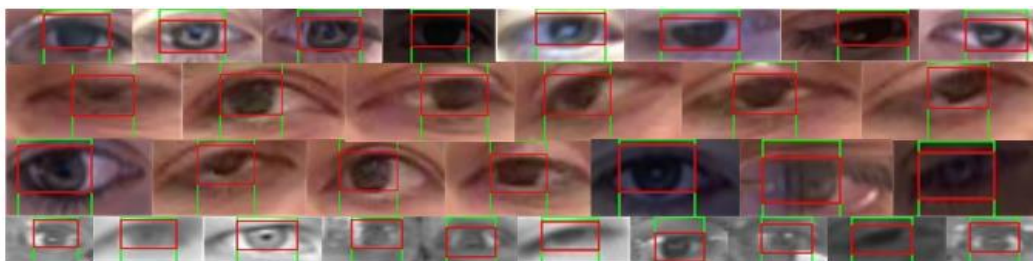441  standard metric similar to S$_{ED}$.

442  It can be analyzed that the model performs comparatively better for Talking-Face and BIO-ID
443  datasets as compared to GI4E dataset based on the corresponding properties (as discussed in Section
444  3). However, there are several crucial aspects to be noted in each case. First, in contrast to *wec*
445  measures in Table 2, the ED($c_a$, $c_e$) error in Table 4 for Talking-Face is 1.96 which is higher than the
446  other two datasets (1.43 and 1.70 for BIO-ID and GI4E respectively) despite the high quality and
447  fewer variations in the former case. This is because the size of images in Talking-Face dataset is
448  comparatively larger than other datasets and consequently, the ED($c_a$, $c_e$) error as well as absolute
449  error ($\mu|x_a-x_e|$, $\mu|y_a-y_e|$) in each axis, are also high. However, results from these metrics (i.e. ED,
450  $\mu|x_a-x_e|$, $\mu|y_a-y_e|$) does not align with results in Table 2 (*wec measure*) and therefore, does not reflect
451  the true measure of the standardized difference between estimated and actual pupil coordinates. In
452  contrast, S$_{ED}$ provides more generic and standard representation of error between the actual and
453  estimated coordinates as a percentage of the eye rectangle size. The S$_{ED}$ error for Talking-Face
454  dataset is 2.49% which is less than 3.98% and 3.87% of BIO-ID and GI4E datasets respectively, and
455  also aligns with the *wec* outcomes in Table 2. As mentioned earlier, S$_{ED}$ represents a standardized
456  distance (i.e. pixels) using current eye-frame without depending upon the second eye or
457  interpupillary distance which is not the case in *wec* measurement. Furthermore, S$_{ED}$ interprets the
458  error in term of pixel distance without using any thresholds (as in case of *wec*) and can be utilized as
459  a standard metric to evaluate the true performance of such models in similar problems.

**Figure 4.** Pupil coordinates estimations (green color) vs actual (red) coordinates within BIO-ID dataset

Figure 4 demonstrates the pupil estimation performance of the proposed model for both left and right eye (*x-axis* and *y-axis*) on BIO-ID dataset. The model indicates a perfect overlapping for both axis and more specifically, at the peak positions which represent the extreme pupil and/or iris positions looking extreme left or right, and top or bottom positions. One of the reasons of such robust overlapping is the use of white paddings in our model that helps the adaptive kernel to achieve maximum overlaps at extreme positions resulting in appropriate matching candidates during horizontal and vertical cascades.



**Figure 5.** Horizontal and vertical convolution-based pupil coordinates localization (in randomly selected images from BIO-ID, GI4E and Talking-Face dataset) for dynamic conditions such as gaze position, eye color, intensity, noise interference, eye size and image resolution

As discussed earlier, a custom kernel might help for optimal representation of iris diversity. Additionally, adaptation of kernel size regarding the eye-frame and dynamic threshold for best candidate selection further improves the reliability of our method specifically in dynamic conditions. Figure 5 demonstrates various test cases of iris/pupil detection using proposed methodology for diverse eye properties and varying environmental conditions (e.g. patterns, gaze direction, varying background, half/full closed eyes, colour, intensity, illuminations, resolution, pupil/iris size, gender, ethnicity etc.). It indicates the robustness of model estimations in both horizontal and vertical convolutions specifically at extreme positions (such as left/right corners, top right, half-closed etc.)

482    Primarily, the proposed method is leveraging the pre-trained Dlib-ml that can locate the facial
483    landmarks efficiently and reliably. It helps to filter out the unnecessary background segments within
484    the input image as well as irrelevant facial components excluding the desired regions that contain
485    exact eye-frames. Secondly, the proposed method uses efficient algorithm to adapt the kernel size in
486    accordance with the eye-frame and padding the eye-frame with white surrounding pixels which
487    further reduce the probability of selecting noisy matched candidates as mentioned by [37, 38]. The
488    use of quantile based dynamic threshold to identify the best matching patch further enhances the
489    reliability in proposed algorithm (e.g. outcomes in Figure 4-5).

490



491    **Figure 6.** The wec measure for different datasets using proposed method

492    Figure 6 shows the performance of the proposed method for pupil coordinate estimation using
493    BIO-ID dataset while varying error threshold, to measure the mean *wec* for both eyes. The
494    visualization indicates accuracy over 90% in all cases (i.e. dataset) while considering the strict
495    constraint of e≤*0.05*. More explicitly, the model indicates that in over 97% of cases with
496    high-resolution images/videos (which are ordinary for current technological advancement), the
497    error in estimated pupil position is less than the diameter of pupil itself. Even in the worst-case
498    scenario (i.e. small-size eye-frames in GI4E dataset), the model achieves above 95% accuracy.



| Ground truth file | #LX | LY | RX | RY |
|---|---|---|---|---|
| For BioID_0000.eye | 232 | 110 | 161 | 110 |

499

500    **Figure 7.** Example of annotation error in BIO-ID dataset

501    It is also imperative to mention that some annotation errors may slightly influence the
502    performance measure even though, this is observed in very few cases. For instance, Figure 7
503    indicates the eye centre coordinates annotations in BIO-ID dataset ($R_x$:161, $R_y$:110) provided by [43]
504    for the right-eye of subject *BioID_0000.eye*. However, the correct values are $R_x$:158, $R_y$:108 (refer to
505    Figure 7) which indicate approximately 2 pixels difference in each axis. This is significant for
506    micro-movements estimation and would affect the model performance substantially (e.g. *wec*, $S_{ED}$).
507    Finally, it can be noted that the proposed model performs initial checks on the current frame
508    quality to assure the existence of exactly two eyes (Algorithm 1) within the identified face rectangle.
509    However, additional constraints can further improve the accuracy specifically, in real-time scenarios
510    and video stream data. For instance, [37] used the DL model to identify the blinking eyes which can

further improve the accuracy of proposed model while filtering out the images/frames without distinctive iris/pupil (i.e. separating the closed eyes not to be analyzed for pupil localization). Additional post-processing constraints such as symmetry constraints over the estimated pupils' coordinates in both eyes might improve the gaze estimation accuracy. This might be useful to improve the eye-state information extraction approaches  such as [7] for the deception detection through facial micro-gestures.

## 5.  Conclusions and Future Works

We proposed a novel pupil estimation method utilising the deep learning based facial landmark detection and an image processing algorithm to determine the eye centre within an image frame. Reliable extraction of the eye-frames within the input image is one of the major advantages of using Dlib-ml. This eliminates most of the background and irrelevant segment of the image which helps to identify the target segment using intelligent image processing. We developed a customized iris kernel using multiple images from various datasets, for its generalized representation. The iris kernel is then convolved with eye-frame in two stages (horizontal and vertical) such that no nested strides are performed by convolution function. White paddings surrounding the kernel as well as eye-frame, proved very helpful for template matching between the kernel and overlapped eye-patches, specifically for the extreme eye positions (e.g. left/right corners). Also, utilising a dynamic threshold for identifying the best-matched patch further contributed to reliability in our method.

Compared to several state-of-the-art pupil detection methods, the proposed approach indicated significant improvements in pupil estimation accuracy specifically, with lower-resolution images and minimum error thresholds. We also introduced a standardized distance metric to measure the relative error in model estimation. This metric can be used regardless of image size and resolution which is not the case with most of the existing validation metrics used in similar works. In future, proposed method will be utilised along with eye-blink detection models, to determine eye gaze, in particular for infraduction iris positions. Our method can be useful in various computer vision applications specifically the one requiring precise pupil and eye centre estimation. For instance, the eye related feature extraction in [7] can be replaced with our method to extract the more reliable and micro-level movements within the eyes to distinguish the truthful and deceptive behaviour. More explicitly, this work is expected to direct several application areas such as human-computer interfaces, gaze estimation, emotion recognition, psychological profiling, fatigue detection, healthcare, visual aid and automated deception detection.

## References

1.  Monforte, P.H.B.; Araujo, G.M.; De Lima, A.A. Evaluation of a New Kernel-Based Classifier in Eye Pupil Detection. Proceeding of 17th IEEE International Conference on Machine Learning and Applications (ICMLA), Orlando, FL, 2018, pp. 380-385.

2.  Al-Rahayfeh, A.; Faezipour, M. Eye tracking and head movement detection: A state-of-art survey. *IEEE J. Transl. Eng. Health Med* 2013, vol. 01, pp. 1-12, doi: 10.1109/JTEHM.2013.2289879.

3.  Guan, X.F.X.; Peli, E.; Liu, H.; Luo, G. Automatic calibration method for driver's head orientation in natural driving environment. *IEEE Trans. Intell. Transp. Syst.*, 2013, vol. 14, no. 01, pp. 303–312.

4.  Horak, K. Fatigue features based on eye tracking for driver inattention system.In Proceeding of 34th Int. Conf. Telecommun. Signal Process. (TSP), Aug. 2011, pp. 593–597.

5. Harischandra, J.; Perera, M.U.S. Intelligent emotion recognition system using brain signals (EEG).Proceeding of IEEE EMBS Conf. Biomed. Eng. Sci. (IECBES), Dec. 2012, pp. 454–459.

6. Hansen, D.W.; Ji, Q. In the eye of the beholder: A survey of models for eyes and gaze. *IEEE Trans. Pattern Anal. Mach. Intell* 2010, vol. 32, no. 03, pp. 478–500.

7. O'Shea, J.; Crockett, K.; Wasiq, K.; Kindynis, P.; Antoniades, A.; Boultadakis, G. Intelligent Deception Detection through Machine Based Interviewing.IEEE International Joint conference on Artificial Neural Networks (IJCNN), Rio de Janeiro, 2018, pp. 1-8, doi: 10.1109/IJCNN.2018.8489392.

8. Waheed, H.; Hassan, S.; Aljohani, N.R.; Hardman, J.; Alelyani, S.; Nawaz, R. Predicting academic performance of students with VLE big data using deep learning models. *Computer in human behavior* 2020, vol. 104, doi: https://doi.org/10.1016/j.chb.2019.106189

9. Fasel, I.R.; Fortenberry, B.; Movellan, J.R. A Generative Framework for Real Time Object Detection and Classification. *Computer Vision and Image Understanding* 2005. vol. 98, no. 01, pp. 182210.

10. Feng, G.C.; Yuen, P.C. Variance Projection Function and Its Application to Eye Detection for Human Face Recognition. *Int'l J. Computer Vision* 1998, vol. 19, pp. 899-906.

11. Feng, G.C.; Yuen, P.C. Multi-Cues Eye Detection on Gray Intensity Image. *Pattern Recognition* 2001, vol. 34, pp. 1033-1046.

12. Kawato, S.; Ohya, J.Two-Step Approach for Real-Time Eye Tracking with a New Filtering Technique. Proceeding of Int'l Conf. System, Man and Cybernetics, 2000, pp. 1366-1371

13. Kawato, S.; Ohya, J. Real-Time Detection of Nodding and Head-Shaking by Directly Detecting and Tracking the BetweenEyes, Proceeding of IEEE Fourth Int'l Conf. Automatic Face and Gesture Recognition, 2000, pp. 40-45.

14. Kawato, S.; Tetsutani, N. Detection and Tracking of Eyes for Gaze-Camera Control. Proceeding of . 15th Int'l Conf. Vision Interface, 2002.

15. Kawato, S.; Tetsutani, N.; Real-Time Detection of Between-theEyes with a Circle Frequency Filter. Proceeding of Asian Conf. Computer Vision '02, 2002, vol. 02, pp. 442-447, 2002.

16. Huang, W.M.; Mariani, R. Face Detection and Precise Eyes Location. Proceeding of Int'l Conf. Pattern Recognition, 2000.

17. Pentland, A.; Moghaddam, B.; Starner, T. View-Based and Modular Eigenspaces for Face Recognition. Proceeding of IEEE Conf. Computer Vision and Pattern Recognition, June 1994.

18. Huang, J.; Wechsler, H. Eye Detection Using Optimal Wavelet Packets and Radial Basis Functions (RBFs). *Int'l J. Pattern Recognition and Artificial Intelligence* 1999, vol. 13, no. 07.

19. Hansen, D.W.; Hansen, J.P.; Nielsen, M.; Johansen, A.S.; Stegmann, M.B. Eye Typing Using Markov and Active Appearance Models. Proceeding of IEEE Workshop Applications on Computer Vision, 2003, pp. 132-136.

20. Hansen, D.W.; Hansen, J.P. Robustifying Eye Interaction. Proceeding of Conference on Vision for Human Computer Interaction, 2006, pp. 152-158.

21. Yiu, Y.H.; Aboulatta, M.; Raiser, T.; Ophey, L.; Flanagin, V.L.; Eulenburg, P.Z.; Ahmadi, S.A. DeepVOG: Open-source pupil segmentation and gaze estimation in neuroscience using deep learning. *Journal of Neuroscience Methods* 2019, vol. 324, 108307.

22. Li, D.; Winfield, D.; Parkhurst, D.J. Starburst: a hybrid algorithm for video-based eye tracking combining feature-based and model-based approaches. IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), 2005, p. 79.

23. Fuhl, W.; Santini, T.; Kasneci, G.; Kasneci, E. Pupilnet: Convolutional Neural Networks for Robust Pupil Detection. 2016,, arXiv:1601.04902 .

24. Villanueva, A.; Ponz, V.; Sanchez, S.; Ariz, M.; Porta, S.; Cabeza, R. Hybrid method based on topography for robust detection of iris centre and eye corners., *ACM Trans. Multimedia Comput. Commun. Appl.* 2013, vol. 04, no. 25, pp. 25:1-20, doi: http://doi.acm.org/10.1145/2501643.2501647.

25. Fuhl, W.; Kübler, T,; Sippel, K.; Rosenstiel, W.; KasneciExcuse, E. Robust pupil detection in real-world scenarios. Azzopardi, G.; Petkov. N (Eds.), *Computer Analysis of Images and Patterns, Springer International Publishing*, 2015, pp. 39-51.

26. Liu, Y.; Hsueh, P.Y.; Lai, J.; Sangin, M.; Nussli, M.A.; Dillenbourg, P. Who is the expert? Analyzing gaze data to predict expertise level in collaborative applications. IEEE International Conference on Multimedia and Expo, 2019, pp. 898-90.

27. Marshall, J.S.P. Identifying cognitive state from eye metrics Aviation", *Space, and Environmental Medicine* 2007, vol. 78, no. 05, pp. 165-175

28. Henderson, J.M.; Shinkareva, S.V.; Wang, J.; Luke, S.G.; Olejarczyk, J. Predicting cognitive state from eye movements. *PLoS One* 2013, vol. 08.

29. Król, M.; KrólA, M.E. Novel approach to studying strategic decisions with eye-tracking and machine learning. *Judgment and Decision Making* 2017, vol. 12, no. 06, pp. 596-609

30. Steichen, B.; Conati, C.; Carenini, G.; Inferring visualization task properties, user performance, and user cognitive abilities from eye gaze data. *ACM Transaction on Interactive Intelligent Systems (TiiS)* 2014, vol. 04, no. 02.

31. Shojaeizadeh, M.; Djamasbi, S.; Paffenroth, R.C.; Trapp, A.C. Detecting task demand via an eye tracking machine learning system. *Decision Support Systems* 2019, vol. 116, pp. 91-101.

32. Krafka, K.; Khosla, A.; Kelnhofer, P.; Kannan, H.; Bhandarkar, S.; Matusik, W.; Torralba, A. Eye tracking for everyone", IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, USA, 2016, doi: 10.1109/CVPR.2016.239.

33. Naqvi, R.A.; Arsalan, M.; Batchuluun, G.; Yoon, H.S.; Park, K.R. Deep Learning-based gaze detection system for automobile drivers using a NIR camera sensor. *Sensors* 2018, vol. 18, no. 02.

34. Vera-Olmos, F.J., Melero, H.; Malpica, N. DeepEye: Deep convolutional network for pupil detection in real environments. *Integrated Computer-Aided Engineering, IOS Press* 2019, vol.26, no. 01, pp. 85-95, doi: 10.3233/ICA-180584.

35. Fuhl, W.; Santini, T.; Kasneci, G.; Kasneci, E. PupilNet: Convolutional Neural Networks for Robust Pupil Detection. ArXiv, abs/1601.04902. 2016, doi: https://arxiv.org/pdf/1601.04902.pdf.

36. Chen, S.; Liu, C. Eye detection using discriminatory Haar features and a new efficient SVM. *Image and Vision Computing* 2015, vol. 33, pp. 68-77, doi: https://doi.org/10.1016/j.imavis.2014.10.007.

37. Borza, D.; Itu, R.; Danescu, R. In the Eye of the Deceiver: Analyzing Eye Movements as a Cue to Deception. *Journal of Imaging, MDPI* 2018, vol. 04, no. 10, pp. 1-20, doi: https://doi.org/10.3390/jimaging4100120.

38. Timm, F.; Barth, E. Accurate eye centre localisation by means of gradients. *Visapp*, 2011, vol. 11, pp. 125–130.

39. King, D.E. Dlib-ml: A machine learning toolkit. *J. Mach. Learn. Res.* 2009, vol. 10, pp. 1755–1758.

40. Valenti, R.; Gevers, T. Accurate eye centre location and tracking using isophote curvature. Proceeding of IEEE, CVPR, Alaska, 2008, pp. 1–8.

41. Cristinacce, D.; Cootes, T.; Scott, I. A Multi-Stage Approach to Facial Feature Detection. Proceeding of British Machine Vision Conf. BMVA Press, 2004, pp. 231-240, doi:10.5244/C.18.30.

42. Cootes, T. Talking Face Video. Available online: http://www-prima.inrialpes.fr/FGnet/data/01-TalkingFace/talking_face.html (accessed on 16 October 2019).

43. Jesorsky, O.; Kirchberg, K.J.; Frischholz, R.W. Robust face detection using the hausdorff distance. Proceeding of International Conference on Audio-and Video-Based Biometric Person Authentication, Halmstad, Sweden, 2001, pp. 90–95.

44. Ponz, V.; Villanueva, A.; Cabeza, R. Dataset for the evaluation of eye detector for gaze estimation. ACM Conf. on Ubiquitous Computing 2012, pp. 681–684.

45. Asadifard, M.; Shanbezadeh, J. Automatic adaptive centre of pupil detection using face detection and cdf analysis. Proceeding of IMECS, Hong Kong, 2010, vol. 01, pp. 130–133.

46. George, A.; Routray, A. Fast and Accurate Algorithm for Eye Localization for Gaze Tracking in Low Resolution Images. In IET Computer Vision, 2016, vol. 10, no. 07, pp.660-669.

47. Fusek R.; Dobeš P. Pupil Localization Using Self-organizing Migrating Algorithm. 2020. In: Zelinka I.; Brandstetter P.; Trong Dao T.; Hoang D.V.; Kim. S. (eds) AETA 2018 - Recent Advances in Electrical Engineering and Related Sciences: Theory and Application. AETA 2018. Lecture Notes in Electrical Engineering, Springer, Cham, vol. 554, pp. 207-216, doi: https://doi.org/10.1007/978-3-030-14907-9_21.