

Highlights

mEBAL2 Database and Benchmark: Image-based Multispectral Eyeblink Detection

Roberto Daza, Aythami Morales, Julian Fierrez, Ruben Tolosana, Ruben Vera-Rodriguez

- We are presenting mEBAL2, the largest eyeblink database.
- Eyeblink detection enhanced through the combination of NIR and RGB images.
- Competitive eyeblink detection, up to 99% accuracy in e-learning.
- Validation in challenging wild environments, achieving state-of-the-art results.



mEBAL2 Database and Benchmark: Image-based Multispectral Eyeblink Detection

Roberto Daza^a, Aythami Morales^a, Julian Fierrez^a, Ruben Tolosana^a, Ruben Vera-Rodriguez^a

^a*Biometrics and Data Pattern Analytics Laboratory, Universidad Autonoma de Madrid, Calle Francisco Tomas y Valiente, 11, Campus de Cantoblanco, Madrid, 28049, Spain*

Article history:

Eyeblink Detection, Eyeblink Database, E-learning, Deep Learning
2000 MSC: 0000, 1111

ABSTRACT

This work introduces a new multispectral database and framework to train and evaluate eyeblink detection in RGB and Near-Infrared (NIR). Our contributed dataset (mEBAL2, multimodal EyeBlink and Attention Level estimation, Version 2) is the largest existing eyeblink database, representing a great opportunity to improve data-driven multispectral approaches for blink detection and related applications (e.g., attention level estimation). mEBAL2 includes 21,100 image sequences from 180 different students (more than 2 million labeled images in total) while conducting a number of e-learning tasks of varying difficulty or taking a real course on HTML initiation through the edX MOOC platform. mEBAL2 uses multiple sensors, including two Near-Infrared (NIR) and one RGB camera to capture facial gestures during the execution of the tasks, as well as an Electroencephalogram (EEG) band to get the cognitive activity of the user and blinking events. Furthermore, this work proposes 3 data-driven approaches as benchmarks for blink detection on mEBAL2, where the architecture based on Convolutional Long Short-Term Memory (ConvLSTM) achieved performances of up to 99%. The experiments explored whether combining RGB and NIR spectrum data improves blink detection in training and architectures that merge both types of data. Experiments showed that the NIR spectrum enhances results, even when only RGB images are available during inference. Finally, the generalization capacity of the proposed eyeblink detectors, along with state-of-the-art eyeblink detection implementations, is validated in wilder and more challenging environments like the HUST-LEBW dataset to show the usefulness of mEBAL2 to train a new generation of data-driven approaches for eyeblink detection.

© 2024 Elsevier Ltd. All rights reserved.

1. Introduction

The act of involuntary closing and opening of the eyelids periodically is defined as eyeblink. The eye is one of the most important organs in the human facial structure for image processing applications from behavior analysis to biometric identification [1, 2]. The eyeblink has proven to be a valuable indicator in various fields such as ocular activity, attention, fatigue, emotions, etc., for this reason, eyeblink detection based on image processing has become essential regarding applications in-

volving human behavior analysis, such as driver fatigue detection [3], attention level estimation [4, 5], dry eye syndrome recovery [6], DeepFakes detection [7], among others.

In the e-learning field, eyeblink detection can be a valuable tool to address certain limitations, particularly, when combined with the latest e-learning platforms [8, 9, 10] that allow collecting students' information to improve security, to guarantee a safe and personalized evaluation, and to adapt dynamically the contents and methodologies to different needs. Eyeblink detection is a useful tool to improve e-learning platforms and get high-quality online education, for at least two reasons. First, since the 70s, studies relate the eyeblink rate with cognitive activity like attention [11, 12]. Recent research suggests that lower eyeblink rates can be associated with high attention peri-

e-mail: roberto.daza@uam.es (Roberto Daza),
aythami.morales@uam.es (Aythami Morales), julian.fierrez@uam.es
(Julian Fierrez), ruben.tolosana@uam.es (Ruben Tolosana),
ruben.vera@uam.es (Ruben Vera-Rodriguez)

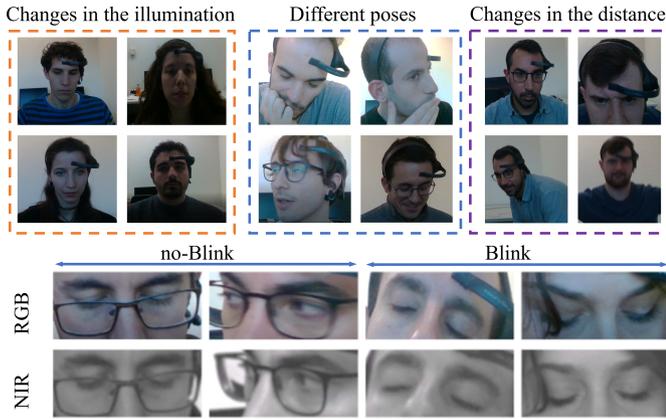


Fig. 1. Different examples from mEBAL2. (top) Sequence images with variations in illumination, posing, and distance to the camera. (bottom) Examples of eyeblink and no-blink with RGB and NIR images.

ods while higher eyeblink rates are related to low attention levels [4, 5, 13]. And second, blink detection can be used in the detection of fraud/cheating/lies and combined with other features like heart rate, gaze tracking, micro-gestures or blood oxygen saturation, can improve the trustworthiness of e-learning platforms.

Over the past few years, significant progress has been made in eyeblink detection [4, 13, 14, 15] thanks to the rise of new computer vision technologies and deep learning techniques, facilitating improved detection of the region of interest and subtle movements of the eyelids, even in challenging conditions such as low lighting and variable poses.

However, the state of the art demonstrates that public eyeblink detectors based on image processing are far from resolving the eyeblink detection problem, HUST-LEBW [14] and MPEblink [15] are recent surveys that demonstrate the need for research in this area. At the moment, there are very few public data-driven algorithms (e.g., neural networks), mainly due to the lack of large eyeblink databases. Most existing datasets useful for research in this area have only a few hundred samples, representing a strong restriction to training data-driven approaches (e.g., deep learning). Also, current public databases restrict their samples to RGB cameras, without using other sensors that have proven to be useful in similar tasks such as NIR cameras in gaze tracking or iris and pupil detection [16].

Considering all of the above, this work aims to provide resources to train and evaluate eyeblink detectors, investigate the potential utility of the NIR spectrum for eyeblink detection, and evaluate new data-driven approaches in realistic environments. The main contributions are:

- We present a new version of the mEBAL database called mEBAL2¹: a multimodal database for eyeblink detection and attention level estimation obtained from an e-learning environment. This database is the largest existing public eyeblink database. This database includes videos from 180 students, with 21,100 labeled image sequences (10,550

eyeblinks and 10,550 no-blink events), more than 2.4 million labeled frames, and students' cognitive activity labels synchronized with all the eyeblink data. Additionally, the new mEBAL2 contains variations on illuminations, poses, distances between user and camera, objects over the face (glasses, hair, hand occlusion, etc.), physical activity, and other naturally-occurring factors. Fig. 1 shows some examples from mEBAL2.

- Specific tasks were designed to provoke changes in the students' cognitive activity like mental load, attention, visual attention, etc. The experiments comprise data from two groups divided into 60 and 120 users. The second group attended a Massive Open Online Course (MOOC) offered by the Universidad Autonoma de Madrid (UAM). This approach provides a real-world environment.
- Our architecture, presented in mEBAL [13], was trained on the proposed mEBAL2 dataset using both RGB and NIR frames. The results demonstrate that NIR images improve the eyeblink detector's training process and outperform the results even when only RGB images are available during the inference.
- Two new data-driven approaches are proposed: An eyeblink detection at frame level based on CNNs that combines all spectral information, and an eyeblink detection at video-sequence level based on ConvLSTM. Additionally, two benchmarks are introduced within mEBAL2: (1) for blink detection at video-sequence level, and (2) for blink detection at frame level.
- Finally, experiments were conducted to showcase the ability of mEBAL2 to train data-driven approaches that can generalize to unseen scenarios. The public benchmark HUST-LEBW, captured in uncontrolled environments, was utilized for evaluation. The results highlight mEBAL2's proficiency in developing new eyeblink detectors that effectively adapt to previously unseen scenarios.

A preliminary version of this article was presented in [13]. This article significantly improves [13] in various aspects:

- An improved version of the mEBAL database with 180 users and 7,550 additional eyeblink events, which now includes 5 times more users and 3.5 times more eyeblink events compared with [13]. Furthermore, a new real e-learning environment has been added.
- For the first time, we trained our architecture proposed in mEBAL on mEBAL2, using RGB, NIR, and both images. We conducted a new study to verify whether using NIR images during training could enhance eyeblink detection.
- Two new blink detectors are proposed, being the most accurate systems on mEBAL2, surpassing the previous version presented in mEBAL.
- Evaluations were conducted on the public HUST-LEBW database [14] using the same eyeblink detector introduced in mEBAL, but with training on the new version. This

¹<https://github.com/BiDALab/mEBAL2>

Table 1. State of the art of eyeblink database.

DB.	Year	Blinks	Users	Att.	Sensors
[17]	NA	61	1	No	RGB
[18]	2007	255	20	No	RGB
[19]	2014	353	4	No	RGB
[20]	2015	300	5	No	RGB
[14]	2019	381	172	No	RGB
[13]	2020	3000	38	Yes	1 RGB - 2 NIR
[15]	2023	8748	NA	No	RGB
Ours	2024	10550	180	Yes	1 RGB - 2 NIR

resulted in an error reduction of 37.31% for the left eye and 27.85% for the right eye.

The rest of the paper is organized as follows. Section 2 summarizes works related to eyeblink detection. Section 3 presents the mEBAL2 database. Section 4 describes the architecture proposal for the eye’s localization and the architectures for the eyeblink verification. Section 5 describes the experiments. Finally, Section 6 provides conclusions and future work.

2. Related Works

2.1. Eyeblink detection databases

There are several well-known databases for blink detection, such as Talking Face [17], ZJU [18], Eyeblink8 [19], and Silesian5 [20], which all share a common limitation: the small number of users and blinks (see Table 1). This limited amount of data implies a significant restriction for data-driven approaches during training. The reported results from published evaluations show saturated scores close to 99% of accuracy due to the reduced number of eyeblink samples [17, 18, 19, 20].

The previous databases are characterized by acquisition setups under controlled environments. The HUST-LEBW database [14] was proposed to explore the detection of eyeblinks in unconstrained scenarios. This database uses eyeblink video clips from 20 commercial movies (The Matrix, Lord of the Rings, etc.). It contains indoor and outdoor examples with scene/illumination changes and varying human poses, similar to a wild environment. HUST-LEBW consists of 172 actors with 381 eyeblinks. It is divided into training with 254 eyeblinks (253 right eye and 243 left eye) and testing with 127 eyeblinks (126 right eye and 122 left eye). The database is unbalanced because the training includes 190 no-eyeblink (190 right eye and 181 left eye) and testing includes 98 no-eyeblink (98 right eye and 98 left eye). The major drawback is the small training set for data-driven approaches (448 samples). A small training set may cause a poor generalization because a few eyeblink examples are not able to handle different head poses, illumination, hair on the eyes, makeup (overstated use of makeup), etc.

The MPEblink database [15] was published recently, and also collected eyeblink data from a wild environment, similar to the HUST-LEBW [14]. This database contains 686 short video clips (7.1–85.9 s) of 86 different movies. 8,748 eyeblink

events were labeled in total, and each video has different eyeblink events from different people.

mEBAL [13] is a previous version of mEBAL2, also captured in an e-learning environment in a multispectral setup (RGB and NIR cameras). mEBAL has 38 users with 3,000 eyeblinks. mEBAL2 comprises three times more users and blink samples with more than 2.4 million frames recorded from 180 MOOC users. Table 1 summarizes the main differences between mEBAL2 and the databases used in the literature.

2.2. Eyeblink detection methods

An eyeblink is a sequence of “eyes open - eyes closed - eyes open” that occurs in a short time. For this reason, eyeblink detectors can be categorized into two groups:

- Eyeblink detection at frame level (image-based eyeblink detection): In this case, the methods classify each frame in open, closed, or the degree of eye closure. Then a sequence of events is defined to detect the eyeblink action like “open/closed/open”, “open/partially-closed/closed/open”, etc. Some methods based on CNN have been proposed for this group. For example, in [21], Anas et al presented two methods based on CNN. The first method had two states (open or closed eyes), and the second had a third state for a partially opened eye. The authors evaluated ZJU and Talking Face datasets getting 93.72% and 100% blink detection accuracy using F1 scores. Phuong et al [22] presented a model based on the Eye Aspect Ratio (EAR), innovating with a custom EAR threshold.
- Eyeblink detection at video-sequence level: This approach uses in a holistic manner the information obtained from a sequence of frames, which is classified entirely as blink or no-blink. Some researchers like Soukupová et al [23] used 13 consecutive frames as input to extract the EAR using facial landmarks. 13 EARs were concatenated as input to an SVM to classify between blink or no-blink. In addition, Hüge et al [14] proposed a model based on a Multi-Scale LSTM (MS-LSTM). Appearance and motion features were extracted in each frame sequence to classify eyeblinks. The MS-LSTM model outperformed some state-of-the-art algorithms in wild environments. Zeng et al [15] presented the InstBlink model recently, which used Query-based methods [24] as its foundation to obtain the face bounding box and eyeblink labels, without using the eye localization method. Zeng et al [25], in a recent work, proposed an approach that captured eyelid movement features to differentiate between blink and no-blink. The network has an architecture with an attention module to generate an attention map, which is fed into a CNN model to jointly learn the appearance and movement features in each frame. Finally, the features extracted are used as input to an LSTM model.

3. mEBAL2 Database

mEBAL2 contains synchronized information from multiple sensors while the students use an interface designed for e-

learning tasks. mEBAL2 acquisition is based on the works of Hernandez et al in [8] and Daza et al in [9]. The authors proposed a platform for remote education assessment called edBB (Biometrics and Behavior for Education). A multimodal acquisition framework was designed to monitor cognitive and eyeblink activity during e-learning tasks. mEBAL2 includes 21,100 events (10,500 blinks and 10,500 no-blinks) from 180 students/sessions. The session duration varies from 15 to 40 min. Each eyeblink event has 19 frames using three cameras: one RGB and two NIR cameras. This database contains 2,405,400 frames (3 cameras \times 19 frames \times 21,100 events \times 2 eyes), making it the largest existing eyeblink database.

Therefore, mEBAL2 provides a dataset consisting of 540 long-duration videos (1 RGB video and 2 NIR videos per session). Each video comes along with the facial bounding box information, 68 facial landmarks, and cropped eye regions for each frame. Furthermore, the dataset includes timestamps for eyeblink and no-eyeblink events and a total of 21,100 cropped samples. Additionally, the dataset provides EEG band information, including attention level, meditation level, 5 electroencephalographic channels, and eyeblink intensity measures. Finally, mEBAL2 contributes two subsets: (i) For frame-level blink detectors, a subset based on the eye state, with 21,000 frames for open eyes and 21,000 frames for closed eyes, is provided. This offers a resource where frame-level blink detectors can be trained and evaluated. (ii) For video-sequence level blink detectors, the subset contains 10,500 blinks and 10,500 no-blinks.

The acquisition setup uses the following sensors (see [9] for a graphical representation of the setup):

- An Intel RealSense (model D435i), which contains 1 RGB and 2 NIR cameras. The 3 cameras are configured to 30 Hz and 1280×720 resolution. According to the Harvard Database of Useful Biological Numbers [26], an average eyeblink ranges between 100 ms–400 ms. Our experience with mEBAL2 reveals that normally eyeblink duration is between 198–263 ms. Therefore, bearing in mind previous studies, our experience, and the setup settings, the average eyeblink can last between 3 to 13 frames.
- An EEG headset by NeuroSky, which measures the power spectrum density of 5 electroencephalographic channels ($\alpha, \beta, \gamma, \delta, \theta$). EEG measures the voltage signals produced usually by synaptic excitations of the dendrites of pyramidal cells in the top layer of the brain cortex [27]. The signals are produced mainly by the number of neurons and fibers fired synchronously [28]. Eyeblinks introduce artifacts that can be easily recognized in EEG signals. In this dataset, the EEG band was used to generate the initial ground truth necessary to label the eyeblink events.

The eyeblinks were labeled using a semi-supervised approach. For that labeling, we first used the eyeblink information provided automatically by the EEG band as candidates for ground truth, and then a human manually checked all the detected events. Without human intervention, the number of eyeblinks detected was 21,484. After the human intervention, the eyeblinks were reduced to 12,032, where 1,482 were labeled

as possible eyeblinks and the remaining eyeblinks were considered ground truth. Each eyeblink in mEBAL2 contains 19 frames in total.

3.1. Tasks

The database was divided into two groups. The first group of 60 students did a series of tasks carefully designed to reach certain goals, and the second group of 120 students did a real lesson from a MOOC entitled “Introduction to Development of Web Applications” (WebApp), available in the edX platform. The lesson is about introduction to HTML coding, where students perform different tasks including watching videos, reading documents, reading and writing HTML code, and performing a final exam.

The tasks for both groups were designed with two goals. First, to generate changes in the students’ cognitive activity such as mental load, attention, visual attention, etc., looking to cause variations of the eyeblink rate. Second, to generate a realistic setting of online assessments. The tasks can be categorized into five groups (see [9] for a video demonstration²):

- Enrollment form: Student’s data are obtained here. This simple task targets a relaxed state with attention levels between normal and low.
- Logical questions: These require more complex interactions, and some of them include crosswords and mathematical problems for the first group. For the MOOC course, some activities involve writing HTML codes and generating more efficient ones.
- Visual tasks: These demand visual attention from the students under different situations, such as watching pre-recorded classes, describing images, and detecting errors in HTML code.
- Reading tasks: Reading documents has proven to have an impact on eyeblink rates and it is highly common in e-learning environments.
- Multiple choice questions: These are essential to help evaluate the students on assessment platforms and most Learning Management Systems provide templates to perform these assessments [29].

4. Proposed Eyeblink Detector

We propose architectures for eyeblink detection: (1) ROI localization, which is commonly used in eyeblink detectors, and (2) Eyeblink verification.

²<https://www.youtube.com/watch?v=JbcL2N4YcDM>

4.1. ROI detection

A sequential approach using deep learning is proposed for ROI detection in 5 steps: (i) Face detection: using the well-known RetinaFace Detector [30], a robust single-stage face detector that uses extra-supervised and self-supervised multi-task learning. (ii) Landmark detection: using a 68 SBR landmark detector [31], based on VGG-16 [32] and Convolutional Pose Machines stages [33]. (iii) Face Alignment: The Dlib library [34] is used to align both eyes parallel to a horizontal line [35]. The alignment is performed utilizing five landmarks (two eyes, the nose, and two at the mouth). The inclination angle formed by the eye landmarks is calculated in the face rotation process, and an affine transformation is applied to correct the tilt. (iv) Data quality: Using detectors' probabilities (p_{ROI}), ROI quality can be assessed. Two different probabilities are calculated, one before alignment ($p_{ROI_{pre}}$) and another after ($p_{ROI_{post}}$), leading to three potential decisions: maintaining alignment ($0.60 \leq p_{ROI_{post}}$), not maintaining it ($0.60 > p_{ROI_{post}}$ or $\frac{2}{3} \times p_{ROI_{pre}} > p_{ROI_{post}}$), or discarding the frame ($0.25 > p_{ROI_{post}}$ and $0.25 > p_{ROI_{pre}}$). Consequently, the process involves recalculating the eye landmarks. A p_{ROI} below 0.25 indicates a failed detection or a turned head. (v) Eye cropping: Finally, each eye was cropped using the ROI's information from the landmark detectors. Later, all eyes were resized to 50×50 .

4.2. Eyeblink verification architecture

The mEBAL2 experimental framework includes 2 blink detectors at frame-level based on CNN architectures and 1 blink detector at video-sequence level based on a ConvLSTM architecture [36].

One-Eye ConvNet architecture (OE-ConvNet) [13]: This architecture is based on the popular VGG16 neural network model [32] (see mEBAL [13] for details). The architecture is formed by 3 convolutional layers with ReLU activation (32/32/64 filters of size 3×3), with 3 max pooling layers between them. The last layer is used as input for a dense layer of 64 units with ReLU activations and 0.5 of dropout. In this work, we have adapted [13] for different training scenarios including RGB and NIR spectrums. During the training process, the RGB and NIR images were introduced in the training batch depending on the scenario.

Left NIR + Right NIR + RGB ConvNet architecture (LI-RI-RGB-ConvNet): We propose a late fusion [37] using the 3 channels of the Intel RealSense information (1 RGB and 2 NIR). It consists of six inputs (2 eyes \times 3 cameras), and each input layer has $50 \times 50 \times C$ dimensions, where C is the number of channels for each used spectrum (3 for RGB and 1 for NIR). All 6 inputs were connected to 6 different convolutional blocks, with the same characteristics of the OE-ConvNet. The outputs of the 6 convolutional blocks were concatenated and connected to a dense layer of 64 units (ReLU activation) and an output layer with one unit (sigmoid activation). A dropout of 0.5 was used.

One-Eye ConvLSTM architecture (OE-ConvLSTM): This architecture processes sequences of 10 input images. It consists of three ConvLSTM layers using recurrent activation

with hard sigmoid, 32 filters with a kernel size of 3×3 (tanh activation) followed by a batch normalization layer and a max pooling layer, excepting the third layer, which changes the number of filters from 32 to 64. Finally, the architecture incorporates a dense layer of 64 units with ReLU activations. Classification between an eyeblink and no-blink is performed by a final output layer with one unit and sigmoid activation. In addition, dropout (0.5) is used.

OE-ConvNet and OE-ConvLSTM offer significant advantages as they were trained to model each eye separately. Therefore, they allow for detecting eyeblinks in side poses, even when one of the eyes is occluded. On the other hand, the LI-RI-RGB-ConvNet takes advantage of all the information provided by the RealSense camera.

5. Experiments and Results

5.1. Experimental protocol

The proposed architectures were trained using mEBAL2 from scratch with a batch size of 32, an Adam optimizer (0.001 learning rate) and a binary cross-entropy loss.

The mEBAL2 benchmark includes a leave-one-out cross validation protocol with one user for testing and the remaining users for training. The process was repeated for each user in the database and the obtained results were averaged. The decision threshold was fixed to the point of Equal Error Rate (EER), in which the False Positive and False Negative rates in blink detection are equal.

The generalization ability of the eyeblink detector trained with mEBAL2 was evaluated on the public benchmark HUST-LEBW [14]. The HUST-LEBW dataset comprises videos obtained from films characterized as in-the-wild completely different to the mEBAL2 environment (e-learning).

5.2. Experiments: mEBAL2 Benchmark

5.2.1. mEBAL2: Blink detection at frame-level

Table 2 presents the mEBAL2 benchmark results of the OE-ConvNet under different training scenarios (e.g., RGB or NIR, different eyes). The results show detection accuracies up to 97.30% in the RGB and 93.94% in the NIR images. Results also show how training a specific detector for each eye leads to a slight improvement compared to training one for both eyes.

One of our goals was to understand if the NIR images could serve to improve current eyeblink detectors. For this reason, we trained OE-ConvNet using RGB, NIR, or both. The results in Table 2 show that eyeblink detection with RGB images is more accurate than the NIR images with a difference of 2.49% approximately. The results with the left NIR camera are similar to the right NIR camera. We trained our OE-ConvNet approach using data from all 3 images (1 RGB + 2 NIR) and evaluated the performance over the RGB images, adopting a similar approach to [38]. The batch of size 32 was generated with both RGB and NIR images. The NIR images (size $50 \times 50 \times 1$) were expanded to size $50 \times 50 \times 3$, which became the input size for our architecture. As we can see in Table 2, there is an improvement of 0.41% when all three images are used for training (FS = Full

Table 2. Comparison of OE-ConvNet on mEBAL2 under different training/evaluation scenarios. The Eyes column denotes the eyes used for training/evaluation. FS: Full Spectrum (RGB and both NIR cameras). NIR_x: x is the Camera (R: Right, L: Left, or B: Both).

Eyes	Training	Evaluation	Acc
Both	RGB	RGB	0.9615
	NIR _R	NIR _R	0.9373
	NIR _L	NIR _L	0.9360
	NIR _B	NIR _B	0.9394
	FS (RGB+NIR _B)	RGB	0.9656
Left	RGB	RGB	0.9730
Right	RGB	RGB	0.9669

Table 3. Eyeblink detection accuracy at frame-level obtained by OE-ConvNet (FS), LI-RI-RGB-ConvNet, and two existing blink detectors [22, 23].

Method	Acc
Blink Detection [22]	0.5837
Soukupova Threshold [23] + Insightface	0.6153
OE-ConvNet (FS)	0.9656
LI-RI-RGB-ConvNet	0.9760

Spectrum), compared to the results obtained when training with RGB only.

Finally, Table 3 presents a comparison between the performance obtained by the OE-ConvNet (FS = Full Spectrum) and LI-RI-RGB-ConvNet EyeBlink detectors (our proposals), and two existing eyeblink detectors: Blink Detection+ [22] and Soukupova Threshold [23] + InsightFace. The thresholds of both detectors [22, 23] were retrained using the frame-level subset of mEBAL2 and the cross-validation protocol proposed in the mEBAL2 benchmark. The face detector of [23] was updated with the state-of-the-art detector InsightFace [30]. Our LI-RI-RGB-ConvNet proposal achieves the highest accuracy with 0.9760, outperforming OE-ConvNet (FS). This demonstrates that our LI-RI-RGB-ConvNet architecture enhances accuracy through late fusion [38] of information from the 3 cameras instead of the heterogeneous training used for OE-ConvNet (FS). However, as a downside, this multispectral approach requires specific hardware with three synchronized sensors. The performance of data-driven approaches such as OE-ConvNet (FS) and LI-RI-RGB-ConvNet is clearly superior to the performance achieved by the methods Blink Detection+ and Soukupova Threshold + InsightFace based on the EAR threshold (see Section 2.2).

The models proposed and evaluated here are aimed to demonstrate the usefulness of mEBAL2 for training and evaluating novel blink detectors, and in that regard compare well with recent methods like [22, 23]. Note that we are not claiming superior performance of our models in comparison with cutting-edge detectors based on advanced data-driven learning architectures such as [15, 25].

Table 4. Eyeblink detection accuracy at video-level obtained by OE-ConvLSTM and two state-of-the-art implementations [22, 23].

Method	Acc
Blink Detection+ [22]	0.6758
Soukupova SVM [23] + Insightface	0.8145
OE-ConvLSTM	0.9909

5.2.2. mEBAL2: Blink detection at video-sequence level

Table 4 presents the eyeblink detection accuracies at video-sequence level of the mEBAL2 Benchmark for our OE-ConvLSTM and two existing blink detectors: Blink Detection+ [22] and Soukupova SVM [23] + InsightFace. Both detectors [22, 23] were retrained on the video-sequence level subset of mEBAL2, including recalibration of Blink Detection+ thresholds and parameter optimization for the Soukupova SVM model. The results demonstrate the potential of mEBAL2 database to train eye detectors obtaining an accuracy of 0.9909. The methods proposed in [22, 23] improve in comparison with the accuracies obtained at frame level. However, the performance of these methods, based on eye landmarks, adaptive thresholds, and SVM is far from those obtained using data-driven learning architectures such as OE-ConvLSTM.

These results are even more remarkable when considering the size of the database and the challenging e-learning environment, where pose changes are common due to the students looking at the keyboard, resulting in closed eyes appearance, as well as changes in lighting, among other factors.

5.3. Experiments on HUST-LEBW: Evaluating the generalization ability of models trained with mEBAL2

Table 5 compares different eyeblink detectors (ours vs relevant related works) based on the public benchmark HUST-LEBW [14]. Our OE-ConvLSTM architecture (Proposal 4) was evaluated through two distinct approaches: (i) OE-ConvLSTM trained using the RGB images of mEBAL2, and (ii) OE-ConvLSTM trained using the HUST-LEBW images. The architecture trained with mEBAL2 achieves the second-best performance in terms of the F1 metric for both eyes, only being outperformed by the recent eyelid method [25]. It is important to note that eyelid incorporates a more complex structure, including an attention generator, CNN, and LSTM architectures (see Section 2.2). When OE-ConvLSTM is trained with mEBAL2, a slight decrease in the performance of approximately 3% in the F1 metric is observed, demonstrating the effectiveness of mEBAL2 for training data-driven approaches capable of generalizing in unseen scenarios.

Our OE-ConvNet architecture (Proposals 1–3) was evaluated with different training settings on mEBAL2: Proposal 1 was trained with both eyes using RGB and NIR images, Proposal 2 was trained with both eyes using only RGB images, and Proposal 3 consists of two detectors trained using RGB images for both eyes separately. Proposal 1 outperforms the best results in the F1 score for architectures trained on mEBAL2, even surpassing our OE-ConvLSTM architecture. These results suggest the usefulness of multispectral training (RGB+NIR) when

Table 5. Eyblink detection results on the HUST-LEBW dataset [14]. Our OE-ConvNet proposals were trained on mEBAL2 (see Table 2 for the training configuration of each Proposal). Our OE-ConvLSTM (Proposal 4) underwent two distinct training: one on mEBAL2 and the other on HUST-LEBW. The method described in [23] was updated using InsightFace [30].

Training	Method	Eye	Recall	Precision	F1
HUST-LEBW	[23]	Both	0.4073	0.8495	0.5506
	[22]	Both	0.5899	0.8005	0.6790
	[15]	Both	0.9764	0.5662	0.7168
	[14]	Left	0.5410	0.8919	0.6735
		Right	0.4444	0.7671	0.5628
	[25]	Left	0.9180	0.8960	0.9069
		Right	0.9127	0.9274	0.9200
	Proposal 4	Left	0.8968	0.8014	0.8464
Right		0.8780	0.7826	0.8276	
mEBAL	[13]	Left	0.9603	0.6080	0.7446
		Right	0.7950	0.7348	0.7637
mEBAL2	Proposal 1	Left	0.9440	0.7564	0.8399
		Right	0.8770	0.7868	0.8295
	Proposal 2	Left	0.9520	0.7126	0.8151
		Right	0.8934	0.6855	0.7758
	Proposal 3	Left	0.9200	0.6928	0.7904
		Right	0.9262	0.7152	0.8072
	Proposal 4	Left	0.8596	0.7656	0.8100
		Right	0.8303	0.7750	0.8017

testing on a different dataset (note that HUST-LEBW includes RGB images only).

Proposal 2 has the same architecture as Proposal 1 but it was trained only with RGB images. As a result, Proposal 2 has lower performance than Proposal 1, especially for the right eye. It is interesting because this indicates that training with more data and with both spectra allows the creation of models with a greater generalization capacity for different environments (different illumination, head orientation, etc.). Furthermore, our initial approach presented in [13] shares the same architecture as Proposal 2 (OE-ConvNet). However, it was trained with the first version of mEBAL with RGB images and therefore presents inferior results in comparison with our Proposal 2 trained now with mEBAL2. Proposal 2 improves the F1 metric in both eyes, with 7.05% for the left eye and 1.21% for the right eye. These results suggest the importance of the usage of wide databases to train data-driven eyblink detectors.

The training of the OE-ConvNet architecture in Proposal 3, which consists of two detectors trained using RGB images for both eyes separately, outperforms the performance of the other OE-ConvNet proposals in the mEBAL2 evaluation (see table 2). However, in the HUST-LEBW evaluation, Proposal 3 obtains worse results for F1 metrics in both eyes than the OE-ConvNet architecture trained with both eyes using RGB and NIR images (Proposal 1). Also, for the left eye, it obtains lower F1 scores than the OE-ConvNet architecture trained with both eyes using

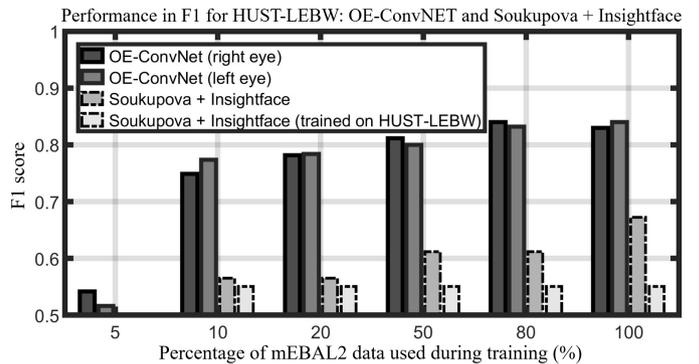


Fig. 2. F1 score results on HUST-LEBW evaluation for different training ratios in mEBAL2 for OE-ConvNet and Soukupova [23] + Insightface architectures.

RGB images (Proposal 2). This result again shows the importance of training with databases with a large number of samples and how the NIR spectrum can be useful to train data-driven approaches with robust generalization capabilities.

Finally, Figure 2 shows the results achieved for the F1 scores on the HUST-LEBW dataset for different training percentages on mEBAL2 of our OE-ConvNet method (Proposal 1) along with our adaptation of Soukupova + InsightFace [23]. The results demonstrate an increase in accuracy for both architectures when training with a larger volume of data. Even Soukupova + InsightFace achieves better performance when trained on mEBAL2 than with HUST-LEBW training. As we can see, the large number of samples and users in mEBAL2 allows for the improvement in the performance of eyblink detectors even in unconstrained scenarios.

6. Conclusions

This work has presented a new multispectral database for eyblink detection. mEBAL2 is 3.52 times wider than the first version in terms of samples and around 5 times larger in terms of users, being the largest existing eyblink database in the literature for research in image- and video-based eyblink detection and related applications, e.g.: attention level estimation [5] and presentation attack detection to face biometrics [39]. mEBAL2 uses visible and NIR spectra (1 RGB and 2 NIR cameras).

Besides, we explored the effects of the visible (RGB) and NIR spectra for eyblink detection. Our results demonstrate that: (i) the approaches trained with both spectra have a good generalization capacity for unseen scenarios, (ii) the combination of the RGB and NIR spectrum through late fusion architectures improves the results in eyblink detection on e-learning environments.

Our proposed architecture for blink detection at video-sequence level, based on ConvLSTM, has achieved the highest levels of accuracy, approximately 99%, in the challenging e-learning environment considered in mEBAL2. Additionally, our methods achieved the second-best performance under uncontrolled conditions in the HUST-LEBW dataset, only surpassed by the eyelid method [25], which is a more complex architecture (see Section 2.2).

mEBAL2 has proven to be a valuable resource to train data-driven algorithms, since a simple CNN learning architecture, when trained on mEBAL2, has demonstrated robust generalization capabilities and significantly improved results compared to its performance with mEBAL. Therefore, the results show that mEBAL2 can be used to train a new generation of data-driven approaches for eyeblink detection.

Future work includes: exploring in more depth the NIR spectrum, advancing in new architectures to leverage the temporal information across frames (GRU, Transformers, etc.), exploiting modern multimodal strategies [40] integrating context information in the periocular region [2], exploiting recent advances in generative face biometrics [41], and exploiting general large-scale AI models with facial analysis capabilities [42] to provide added-value in this specific problem of eyeblink detection.

7. Acknowledgments

Support by project HumanCAIC (TED2021-131787B-I00 MICINN), BIO-PROCTORING (GNOSS Program, Agreement Ministerio de Defensa-UAM-FUAM dated 29-03-2022) and Catedra ENIA UAM-VERIDAS en IA Responsable (NextGenerationEU PRTR TSI-100927-2023-2). Research partially funded by the Autonomous Community of Madrid. Roberto Daza is supported by a FPI fellowship from MINECO/FEDER. A. Morales is supported by the Madrid Government (Comunidad de Madrid-Spain) under the Multiannual Agreement with Universidad Autónoma de Madrid in the line of Excellence for the University Teaching Staff in the context of the V PRICIT (Regional Programme of Research and Technological Innovation).

References

- [1] F. Alonso-Fernandez, R. A. Farrugia, J. Fierrez, J. Bigun, Super-resolution for selfie biometrics: Introduction and application to face and iris, in: A. Rattani, R. Derakhshani, A. Ross (Eds.), *Selfie Biometrics*, Springer, 2019, pp. 105–128.
- [2] F. Alonso-Fernandez, J. Bigun, J. Fierrez, N. Damer, H. Proença, A. Ross, Periocular biometrics: A modality for unconstrained scenarios, *IEEE Comput.* (2024).
- [3] L. M. Bergasa, J. Nuevo, M. A. Sotelo, R. Barea, M. E. Lopez, Real-time system for monitoring driver vigilance, *IEEE Transactions on Intelligent Transportation Systems* 7 (1) (2016) 63–77.
- [4] R. Daza, D. DeAlcala, A. Morales, R. Tolosana, R. Cobos, J. Fierrez, ALEBk: Feasibility study of attention level estimation via blink detection applied to e-learning, in: *Proc. AAAI Workshop on Artificial Intelligence for Education*, 2022.
- [5] R. Daza, L. F. Gomez, A. Morales, J. Fierrez, R. Tolosana, R. Cobos, J. Ortega-Garcia, MATT: Multimodal Attention Level Estimation for e-learning Platforms, in: *Proc. AAAI Workshop on Artificial Intelligence for Education*, 2023.
- [6] M. Rosenfield, Computer vision syndrome: a review of ocular causes and potential treatments, *Ophthalmic and Physiological Optics* 31 (5) (2011) 502–515.
- [7] T. Jung, S. Kim, K. Kim, Deepvision: Deepfakes detection using human eye blinking pattern, *IEEE Access* (2020) 83144–83154.
- [8] J. Hernandez-Ortega, R. Daza, A. Morales, J. Fierrez, J. Ortega-Garcia, edBB: Biometrics and Behavior for Assessing Remote Education, in: *Proc. AAAI Workshop on Artificial Intelligence for Education*, 2020.
- [9] R. Daza, A. Morales, R. Tolosana, L. F. Gomez, J. Fierrez, J. Ortega-Garcia, edBB-Demo: Biometrics and Behavior Analysis for Online Educational Platforms, in: *Proc. AAAI Conf. on Artificial Intelligence (Demonstration)*, 2023.
- [10] Á. Becerra, R. Daza, R. Cobos, A. Morales, M. Cukurova, J. Fierrez, M2LADS: A System for Generating MultiModal Learning Analytics Dashboards in Open Education, in: *Proc. Annual Computers, Software, and Applications Conference (COMPSAC) in the Workshop on Open Education Resources*, 2023.
- [11] J. Bagley, L. Manelis, Effect of awareness on an indicator of cognitive load, *Perceptual and Motor Skills* 49 (2) (1979) 591–594.
- [12] M. K. Holland, G. Tarlow, Blinking and mental load, *Psychological Reports* 31 (1) (1972) 119–127.
- [13] R. Daza, A. Morales, J. Fierrez, R. Tolosana, mEBAL: A Multimodal Database for Eye Blink Detection and Attention Level Estimation, in: *International Conference on Multimodal Interaction*, 2020, pp. 32–36.
- [14] G. Hu, Y. Xiao, Z. Cao, L. Meng, Z. Fang, J. T. Zhou, J. Yuan, Towards real-time eyeblink detection in the wild: Dataset, theory and practices, *IEEE Transactions on Information Forensics and Security* (2019) 2194–2208.
- [15] W. Zeng, Y. Xiao, S. Wei, J. Gan, X. Zhang, Z. Cao, Z. Fang, J. T. Zhou, Real-time multi-person eyeblink detection in the wild for untrimmed video, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13854–13863.
- [16] F. Alonso-Fernandez, K. Raja, R. Raghavendra, C. Busch, J. Bigun, R. Vera-Rodriguez, J. Fierrez, Cross-sensor periocular biometrics for partial face recognition in a global pandemic: Comparative benchmark and novel multialgorithmic approach, *Information Fusion* 83-84 (2022) 110–130.
- [17] Talking Face, Talking face, https://personalpages.manchester.ac.uk/staff/timothy.f.cootes/data/talking_face/talking_face.html, accessed: 2023-05-25 (2021).
- [18] G. Pan, L. Sun, Z. Wu, S. Lao, Eyeblink-based Anti-spoofing in Face Recognition from a Generic Webcam, in: *Proc. IEEE International Conference on Computer Vision*, 2007.
- [19] T. Drutarovsky, A. Fogelton, Eye Blink Detection using Variance of Motion Vectors, in: *Proc. European Conference on Computer Vision*, 2014, pp. 436–448.
- [20] K. Radlak, M. Bozek, B. Smolka, Silesian Deception Database: Presentation and Analysis, in: *Proc. ACM on Workshop on Multimodal Deception Detection*, 2015, pp. 29–35.
- [21] E. R. Anas, P. Henriquez, B. J. Matuszewski, Online Eye Status Detection in the Wild with Convolutional Neural Networks, in: *Proc. International Conf. on Computer Vision Theory and Applications*, 2017, pp. 88–95.
- [22] T. T. Phuong, L. T. Hien, N. D. Vinh, An Eye Blink Detection Technique in Video Surveillance based on Eye Aspect Ratio, in: *Proc. Advanced Communication Technology*, 2022, pp. 534–538.
- [23] T. Soukupová, J. Cech, Eye Blink Detection using Facial Landmarks, in: *Proc. Computer Vision Winter Workshop*, 2016.
- [24] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end Object Detection with Transformers, in: *Proc. European Conference on Computer Vision (ECCV)*, 2020.
- [25] W. Zeng, Y. Xiao, G. Hu, Z. Cao, S. Wei, Z. Fang, J. T. Zhou, J. Yuan, Eyelid’s intrinsic motion-aware feature learning for real-time eyeblink detection in the wild, *Transactions on Information forensics and security* 18 (2023) 5109–5121.
- [26] H. Richard Schiffman (Ed.), *Sensation and Perception: An Integrated Approach*, John Wiley & Sons, 1990.
- [27] T. Kirschstein, R. Köhling, What is the source of the EEG?, *Clinical EEG and Neuroscience* 40 (3) (2009) 146–149.
- [28] J. E. Hall, M. E. Hall (Eds.), *Guyton and Hall Textbook of Medical Physiology e-Book*, Elsevier Health Sciences, 2020.
- [29] T. Govindasamy, Successful implementation of e-learning: Pedagogical considerations, *The Internet and Higher Education* 4 (3-4) (2001) 287–299.
- [30] J. Deng, et.al, RetinaFace: Single-Shot Multi-Level Face Localisation in the Wild, in: *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2020, pp. 5203–5212.
- [31] X. Dong, S.-I. Yu, X. Weng, S.-E. Wei, Y. Yang, Y. Sheikh, Supervision-by-Registration: An Unsupervised Approach to Improve the Precision of Facial Landmark Detectors, in: *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 360–368.

- [32] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint (2014).
- [33] S.-E. Wei, V. Ramakrishna, T. Kanade, Y. Sheikh, Convolutional Pose Machines, in: Proc. IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 4724–4732.
- [34] Dlib, Dlib c++ library, <http://dlib.net/>, accessed: 2024-02-20 (2024).
- [35] P. Tome, R. Vera-Rodriguez, J. Fierrez, J. Ortega-Garcia, Facial soft biometric features for forensic face recognition, *Forensic Science International* 257 (2015) 171–284.
- [36] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, W.-c. Woo, Convolutional lstm network: A machine learning approach for precipitation nowcasting, *Advances in Neural Information Processing Systems* 28 (2015).
- [37] J. Fierrez, A. Morales, R. Vera-Rodriguez, D. Camacho, Multiple classifiers in biometrics. part 1: Fundamentals and review, *Information Fusion* 44 (2018) 57–64.
- [38] K. Fu, D.-P. Fan, G.-P. Ji, Q. Zhao, J. Shen, C. Zhu, Siamese network for rgb-d salient object detection and beyond, *Transactions on pattern analysis and machine intelligence* 44 (9) (2021) 5541–5559.
- [39] J. Hernandez-Ortega, J. Fierrez, A. Morales, J. Galbally, Introduction to presentation attack detection in face biometrics and recent advances, in: S. Marcel, J. Fierrez, N. Evans (Eds.), *Handbook of Biometric Anti-Spoofing*, Springer, 2023, 3rd Ed.
- [40] A. Peña, I. Serna, A. Morales, J. Fierrez, A. Ortega, A. Herrarte, M. Alcantara, J. Ortega-Garcia, Human-centric multimodal machine learning: Recent advances and testbed on ai-based recruitment, *SN Computer Science* 4 (5) (2023) 434.
- [41] P. Melzi, R. Tolosana, R. Vera-Rodriguez, M. Kim, C. Rathgeb, X. Liu, I. DeAndres-Tame, A. Morales, J. Fierrez, et al., FRCSyn-onGoing: Benchmarking and comprehensive evaluation of real and synthetic data to improve face recognition systems, *Inf. Fusion* 107 (2024) 102322.
- [42] I. Deandres-Tame, R. Tolosana, R. Vera-Rodriguez, A. Morales, J. Fierrez, J. Ortega-Garcia, How good is ChatGPT at face biometrics? a first look into recognition, soft biometrics, and explainability, *IEEE Access* 12 (2024) 34390–34401.