

A Weakly-Supervised Framework for Interpretable Diabetic Retinopathy Detection on Retinal Images

Pedro Costa*, Adrian Galdran, Asim Smailagic, and Aurélio Campilho

Abstract—Diabetic Retinopathy (DR) detection is a critical retinal image analysis task in the context of early blindness prevention. Unfortunately, in order to train a model to accurately detect DR based on the presence of different retinal lesions, typically a dataset with medical expert's annotations at the pixel level is needed. In this paper, a new methodology based on the Multiple Instance Learning (MIL) is developed in order to overcome this necessity by leveraging the implicit information present on annotations made at the image level. Contrary to previous MIL-based DR detection systems, the main contribution is that the proposed technique jointly optimizes the instance encoding and the image classification stages. In this way, more useful mid-level representations of pathological images can be obtained. The explainability of the model decisions is further enhanced by means of a new loss function enforcing appropriate instance and mid-level representations. A comprehensive quantitative and qualitative experimental evaluation on several publicly available datasets is provided, confirming that the proposed technique achieves results comparable or better than other recently proposed methods, while improving the interpretability of the produced decisions.

Index Terms—Multiple Instance Learning, Diabetic Retinopathy Detection, Bag of Visual Words.

I. INTRODUCTION

THE retina is a well-known source of biomarkers that enable the early identification of several human disorders, such as hypertension, heart diseases, or Diabetic Retinopathy (DR), among others. DR is known to be the leading cause of preventable blindness, affecting more than 415 million people worldwide [1]. Fortunately, DR can be detected at its early stages by expert ophthalmologists through routine analysis of the eye fundus [2]. Timely DR detection can lead to the administration of preventive treatments and efficient therapies to avoid vision impairment and further consequences.

To provide early disease diagnosis and appropriate eye care, large-scale global screening programs are often implemented by hospitals and local authorities [3], [4] with great success. In the context of such programs, patients are called to clinical settings in order to acquire eye color images with a retinal fundus camera. These images are then submitted to specialists, who look for visual signs of the presence of lesions, and

P. Costa, A. Galdran, and A. Campilho are with INESC TEC, R. Dr. Roberto Frias, 4200-465, Porto, Portugal, e-mails: {pvcosta, adrian.galdran}@inesctec.pt.

A. Smailagic is with Institute for Complex Engineered Systems, Carnegie Mellon University, Pittsburgh, PA, United States, e-mail: asim@cs.cmu.edu.

A. Campilho is also with Faculdade de Engenharia, Universidade do Porto, R. Dr. Roberto Frias, 4200-464, Porto, Portugal, e-mail: campilho@fe.up.pt.

* Corresponding Author.

Manuscript received July 31, 2017; revised July 31, 2017.

perform diagnosis based on this. A sample of these potential signs of disease is shown in Fig. 1.

Unfortunately, more than 83% of undiagnosed DR patients are located in underdeveloped areas, where there is a lack of specialists to attend large masses of population, blocking the appropriate implementation of screening programs [5]. For this reason, Computer-Aided Diagnosis (CAD) systems capable of detecting signs of DR from standard retinal fundus images are becoming highly relevant in recent years [6]. An effective automatic DR detection system can substantially reduce the workload experimented by ophthalmologists in the context of large-scale screening programs, having a large positive impact on population healthcare [7], [8].

However, in order to design a CAD system to support expert's decisions in an appropriate manner, such a system must enjoy several desirable properties. First, the amount of annotated data needed to train it must be medium-to-moderate, since manually labeling each pixel on image regions containing lesions can be a time-consuming and error-prone process. Ideally, a CAD system must be able to learn to detect disease signs out of a set of images labeled with a single number indicating the presence or not of DR. Labeling retinal images in this way is a much easier and faster to accomplish task for human experts, and there is already a large quantity of visual data stored at hospitals that has been annotated with this kind of labels, representing an immense source of training data. We refer to this kind of annotations as *weakly labeled* data. Second, any CAD system supporting ophthalmologists' clinical decisions must work reliably and in an interpretable manner, in order to fit regular clinical work-flows.

To address both of these challenges, the main contribution of this paper consists of a new technique for DR detection capable of learning from a set of weakly labeled images and performing interpretable diagnose prediction. The proposed technique allows to train a DR detection CAD system at an image level using implicit local information, *e.g.* decide if the image is healthy or not based on lesions present in certain image regions, even if their location is unknown. This is achieved by means of a novel Multiple Instance Learning (MIL) technique that improves upon previously proposed MIL approaches by jointly learning to encode and classify visual information coming from localized areas of the image. As a second contribution, the interpretability of the proposed model is enhanced by means of a constraint imposed on the learned representations, that forces them to remain sparse in case of healthy images, while becoming dense whenever the image has DR signs. Hence, the proposed system does not require strong manual annotations to be trained, but it can still signal to the

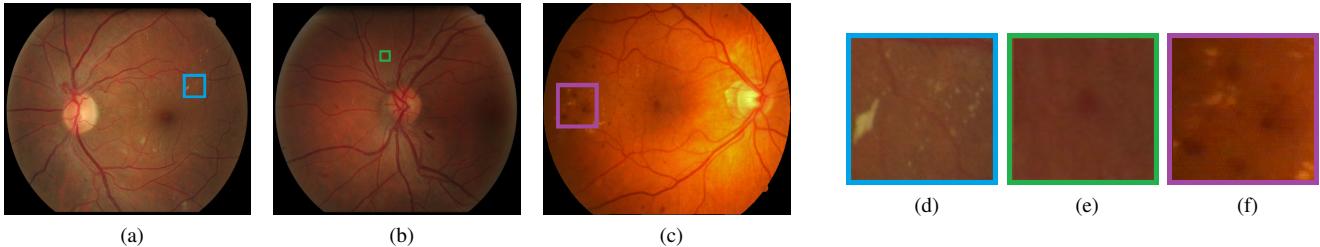


Fig. 1. (a–c) A Retinal Images showing signs of DR (b–d) Lesions associated to DR (d) Exudates (e) Microaneurysms (f) Hemorrhages

90 regions on the image that triggered the diagnosis decision, resulting in a highly interpretable CAD system. Comprehensive
91 performance evaluation on several publicly available datasets
92 favor the proposed technique, demonstrating that it competes
93 well with other recent approaches while providing an increased
94 interpretability outcome. The method presented in this paper
95 is a substantial extension of the conference [publication](#) [9].

II. RELATED WORK

A. DR Detection on Retinal Images

99 DR detection on images of the eye fundus is usually
100 achieved by first locating specific disease signs and lesions
101 in the retinal images. This is typically accomplished based on
102 a conventional machine learning pipeline for detecting objects
103 of interest within images. Given a dataset of image regions
104 containing manually delineated lesions:

- 105 1) **Lesion Description:** Visual features are extracted to
106 characterize each type of lesion modeling their geometric,
107 textural and color appearance.
- 108 2) **Classifier Training:** A classifier is trained to distinguish
109 lesions based on the extracted features.
- 110 3) **Lesion Candidates Extraction:** Given a new image,
111 candidate regions are extracted from it, and the retrieved
112 candidates are described with those same features.
- 113 4) **Lesion Candidate Classification:** These descriptors are
114 inputted to the classifier, which decides first if the candidate
115 is a lesion or not (false positive removal) and/or the
116 most likely type of lesion.

117 Usually, these techniques are specifically designed to deal
118 with a single type of lesion, *e.g.* microaneurysms [10] or
119 hard exudates detection [11]. More generally, red lesions [12],
120 [13], or bright lesions [14], [15], can be detected. After lesion
121 detection has been performed, DR detection and grading can
122 be realized. Several [papers](#) have thus proposed DR detection
123 techniques consisting of combining distinct lesion detection
124 techniques to extract all relevant anomalies, and then merge
125 the results into a single **outcome indicating** the presence or
126 severity of DR [16], [17].

127 The main drawback of the above approach is that it requires
128 an image database that has been previously annotated by
129 a specialist at the lesion level. The lesion borders need to
130 be marked pixel by pixel or with specialized visual tracing
131 tools. This represents a tedious and time-consuming process.
132 Yet another relevant limitation of lesion-detection based DR

133 detection is the necessity of complex and often error-prone pre-
134 processing techniques. For instance, the optic disc typically
135 needs to be located and removed in order to avoid the
136 generation of candidate regions that can be confused as bright
137 lesions [16].

138 To avoid the need for manually segmented training exam-
139 ples, this work differs from multi-lesion detection approaches
140 by framing the problem within the MIL paradigm. MIL allows
141 to build a learning system on which examples can be weakly
142 labeled: only a single indicator is required for a given image,
143 but predictions are formulated based implicitly on region-
144 level characteristics of it. Below we provide a brief theoretical
145 introduction to the MIL framework.

B. Multiple-Instance Learning

146 The MIL framework for binary classification problems
147 considers two main entities, called *bags* and *instances*. In
148 this setting, a bag is composed of an undetermined number
149 of instances. While the goal of a MIL algorithm is to classify
150 bags into positive or negative, it is assumed that instances carry
151 useful information regarding bags containing them. However,
152 the only available ground-truth in MIL is associated to bags.
153 The goal in this approach then becomes to model the implicit
154 relationship between instances and their corresponding bags.
155 A typical example of such a relationship would be: if a bag
156 contains at least a positive instance, regardless of how many
157 negative instances it may contain, it should be declared as
158 positive, whereas a bag should be predicted as negative if it
159 contains only negative instances.

160 There are two main MIL algorithms categories, namely
161 instance-level techniques (ILT) and bag-level techniques
162 (BLT). In ILT, a classifier is trained to classify instances, and
163 instance-level predictions are aggregated to build a bag-level
164 prediction. Examples of this approach are mi-SVM [18], or
165 MIL-Boost [19]. The way in which instance-level predictions
166 are combined will determine the instance/bag relationship
167 being modeled. Following the above example, instance-level
168 predictions can be aggregated with a max-rule: the bag-level
169 prediction is given by the top positive instances contained on
170 it. The main disadvantage of this approach is that not always
171 a single instance should condition the bag-level prediction, *i.e.*
172 a larger set of instances may influence it.

173 BLT differ from ILT in that the classifier is not trained to
174 classify instances, but rather it learns to classify bags directly.
175 The main difference lies in the moment the instance-level

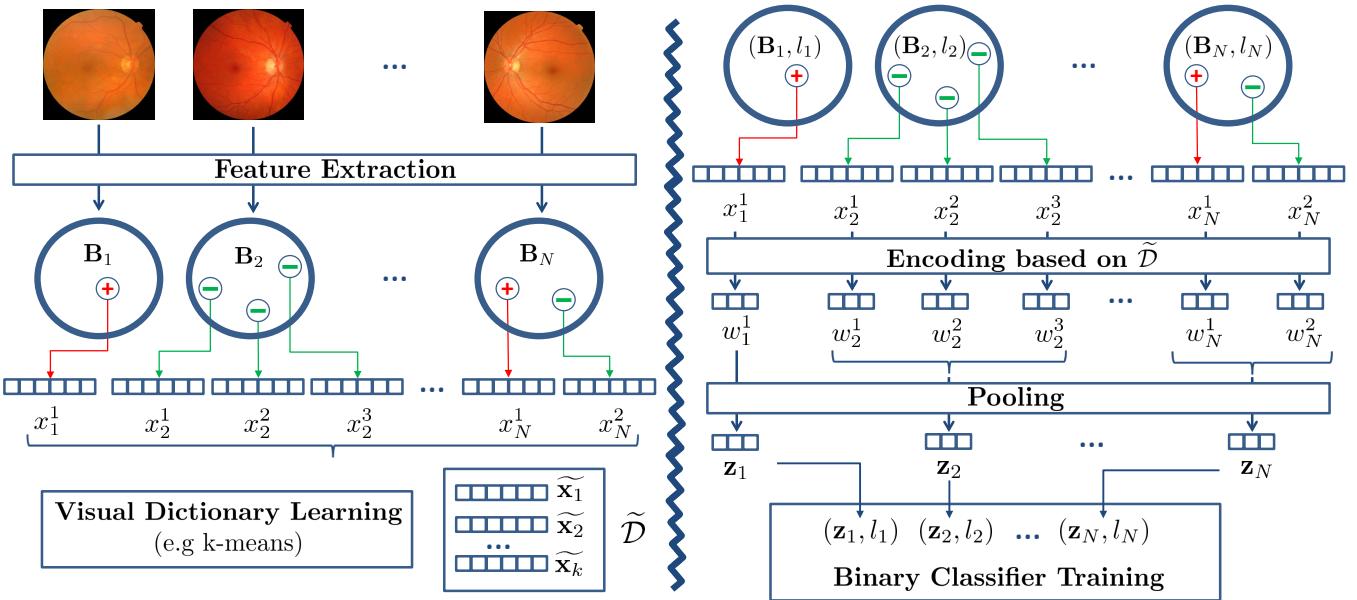


Fig. 2. Conventional BoVW framework. Note that the visual dictionary is learned in a separate stage. Hence, the learned mid-level representations \mathbf{z}_i are completely independent of the binary classifier's performance.

information is aggregated. While ILT combine instance-level predictions, in BLT a bag-level representation is built out of a combination of instance-level representations, and the classifier is trained on this combined representation [20].

The main problem of MIL-BLT is that the amount of instances within a bag is not known a priori, which gives rise to bag representations with varying dimensions. To overcome this obstacle, typically all the instance-level representations of different bags are mapped into a common space. This is achieved with embedding functions followed by pooling operations, and the goal becomes finding a representation space as discriminative as possible.

One particularly interesting MIL-based image classification model is the Bag of Visual Words (BoVW), introduced in [21] for video retrieval. In BoVW, an image (bag) is decomposed into a set of local low-level visual descriptors (instances). These are then mapped onto a common representation, defined by a visual dictionary.

In BoVW techniques, the way the visual dictionary is learned is a critical part of the method. The most popular approaches involve applying unsupervised techniques, such as k -means clustering, on features extracted from a group of images. In this case, the resulting k centroids conform the visual words composing the dictionary. Once the instances associated to an image have been encoded using the visual dictionary, they are combined together via a pooling operation, resulting in a feature vector that is supplied to a standard classifier.

In summary, BoVW is characterized by two separate stages. The first one extracts features from all images and learns a visual dictionary. The second one is composed of four processes: 1) Feature extraction to build instance representations. 2) Encoding of instance representations into a discriminative space. 3) Pooling the encoded representations into a mid-level

representation for each image, and 4) Classifier training on these mid-level representations. An illustration of this process is shown in Fig. 2.

MIL techniques, and in particular BoVW, have been previously proposed with success for medical imaging applications [22]. For instance, MIL was applied in [23] for obstructive pulmonary disease detection on lung CT scans, in [24] for segmentation and diagnosis of histopathology images, or for detecting early signs of dementia on brain MRI in [25]. MIL has also been proposed for retinal image analysis tasks [26], [27]. In this context, the closer techniques to the method proposed in this paper are [28] and [29]. In [28], a MIL-based DR system is proposed, based on a complex pipeline involving multi-scale patch extraction and alternate local-global weight updating to optimize distances between relevant instances in the feature space. In [29], the authors introduce a BoVW technique for DR detection based on sparse SURF features with a semi-soft encoding scheme and max-pooling.

The approach proposed in this paper is also based on the BoVW framework. However, we depart from the conventional two-stage process which: firstly, learns a visual dictionary, and secondly, trains a binary classifier on mid-level bag representations. This is achieved by simultaneously learning to encode the instance-level feature vectors onto useful representations and learning to classify bags based on them. Thanks to this joint learning process, the learned representations are enforced to be useful for the classification task. This way, the classification performance directly drives the learning process from end to end. Furthermore, the mid-level representations are also constrained via a new loss function that is designed to enhance their interpretability. An overview of the proposed system is illustrated on Fig. 3.

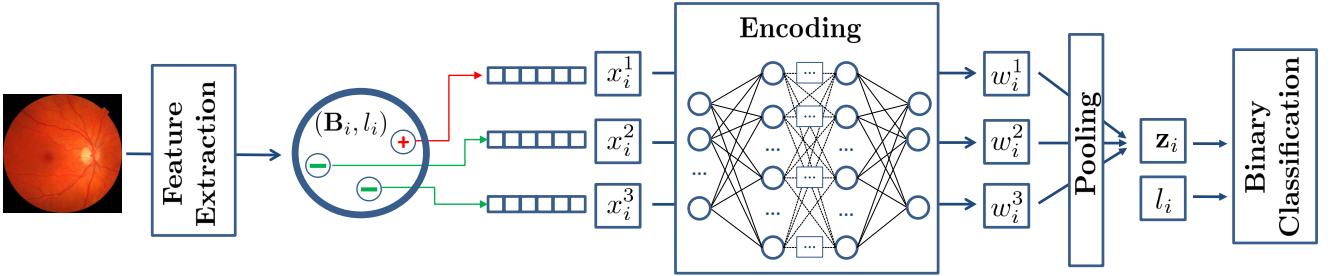


Fig. 3. Improved BoVW algorithm. The system avoids the explicit creation of a visual dictionary. Furthermore, note that in the proposed version of the BoVW the classifier's performance drives the selection of optimal mid-level representations \mathbf{z}_i , as opposed to the conventional approach.

III. BOVW FOR INTERPRETABLE DR DETECTION

To formalize the BoVW approach for MIL in the context of DR detection, let us consider a training dataset $\{(\mathbf{B}_n, l_n)\}_{n=1}^N$ of N retinal eye fundus images \mathbf{B}_n with associated labels l_n , indicating whether they contain pathological signs. From each image \mathbf{B} , $N(\mathbf{B})$ instances are extracted, consisting of a set of local descriptors from a variable number of image regions.

In this case, each bag \mathbf{B} is modeled as follows:

$$\mathbf{B} \approx \{\mathbf{x}^i, 1 \leq i \leq N(\mathbf{B})\}, \quad (1)$$

where $\mathbf{x}^i \in \mathbb{R}^d$ is a feature vector describing the i -th instance found in \mathbf{B} .

In order to classify a new bag, we need to train a binary classifier. However, since $N(\mathbf{B})$ varies for each image \mathbf{B} , the description in eq. (1) is not suitable for this task. The conventional BoVW approach proceeds by extracting the representations for all images in a training set and aggregating them, obtaining a set \mathcal{D} of $(\sum_{i=1}^N N(\mathbf{B}_i))$ d -dimensional descriptors. On this set, an unsupervised clustering technique can be applied, e.g. k -means. In this case, the set of descriptors is summarized into k centroids, known as visual words, which compose the visual dictionary $\tilde{\mathcal{D}} = \{\tilde{\mathbf{x}}^1, \tilde{\mathbf{x}}^2, \dots, \tilde{\mathbf{x}}^k\}$.

Once a suitable dictionary $\tilde{\mathcal{D}}$ is learned, then for every training bag \mathbf{B} the method encodes each of its instances \mathbf{x}^i into a set of k -dimensional codes \mathbf{w}^i based on $\tilde{\mathcal{D}}$ by means of an embedding function $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$. All the $N(\mathbf{B})$ codes of different instances extracted from \mathbf{B} are finally pooled in order to obtain a single k -dimensional representation \mathbf{z} . This can be achieved for instance with the max-pooling operation $P : \mathbb{R}^{N(\mathbf{B})} \times \mathbb{R}^k \rightarrow \mathbb{R}^k$. In this case, the m -th element in \mathbf{z} is given by:

$$z_m = \max_{1 \leq j \leq N(\mathbf{B})} w_m^j. \quad (2)$$

We refer to this final pooled vector \mathbf{z} as the bag's mid-level representation associated with image \mathbf{B} . After processing the training set, the resulting set of mid-level representations, together with the corresponding bag-level labels, are finally supplied to a classifier $\mathcal{C}(\mathbf{B})$, which is trained to discriminate between positive or negative bags.

A. Jointly Learning to Encode Instances and Classify Bags

The standard BoVW strategy outlined above presents several important disadvantages. A relevant deficiency is that the

visual dictionary construction, which determines the resulting mid-level representations, is a process completely isolated from the training of the binary classifier. To compensate for this and build a sufficiently expressive dictionary that can capture the complexity of the feature space, its size is typically large. In some cases the amount of considered visual words can reach thousands [26]. Existing alternatives comprise for instance considering a hierarchical coarse-to-fine dictionary learning, on which an initial fine-grained dictionary is iteratively refined employing the bag-level labels, until an optimal size is reached [30]. However, this is a cumbersome step requiring itself a separate optimization process.

To overcome these drawbacks, a new strategy to simultaneously learn the encoding and classification steps is proposed in this work. In order to avoid the construction of a visual dictionary, two neural networks are built: the first one, $\mathbf{U}(\mathbf{x}; \theta_U)$, learns optimal weights θ_U to produce useful mid-level representations \mathbf{z} while the second, $\mathbf{D}(\mathbf{z}; \theta_D)$ receives those representations and its parameters θ_D are optimized to perform accurate bag-level classification. The error of \mathbf{D} is back-propagated directly to \mathbf{U} , influencing the way in which the mid-level representations are produced by it. In this way, both processes can benefit from each other. An overview of this improved BoVW approach is shown in Fig. 3.

Technically, the neural network $\mathbf{U}(\mathbf{x}; \theta)$ is defined by a series of layers $j \in \{1, \dots, L\}$ with M_j hidden neurons that perform simple linear operations specified by weights θ^j on their inputs \mathbf{x} , followed by a non-linear operation:

$$\mathbf{x} \mapsto \sigma(\theta^j \cdot \mathbf{x}) \quad (3)$$

where the first element of \mathbf{x} is set as $x_0 = 1$ in such a way that θ_0^j contains the bias term. In the above equation, σ denotes a sigmoid function or any other kind of non-linearity.

For a given bag \mathbf{B} , each of its instances \mathbf{x}^i going through \mathbf{U} will be encoded into a M_L -dimensional code vector \mathbf{w}^i . Note that the last layer of \mathbf{U} is followed by a softmax activation function, ensuring that $\sum_{j=1}^k w_j^i = 1$. The set of codes $\{\mathbf{w}^i, 1 \leq i \leq N(\mathbf{B})\}$ computed from every instance in \mathbf{B} are then pooled into the mid-level representation \mathbf{z} .

Regarding the pooling stage, several options can be applied, such as average pooling, which averages all codes in $\{\mathbf{w}^i, 1 \leq i \leq N(\mathbf{B})\}$. However, this can lead to a smoothing

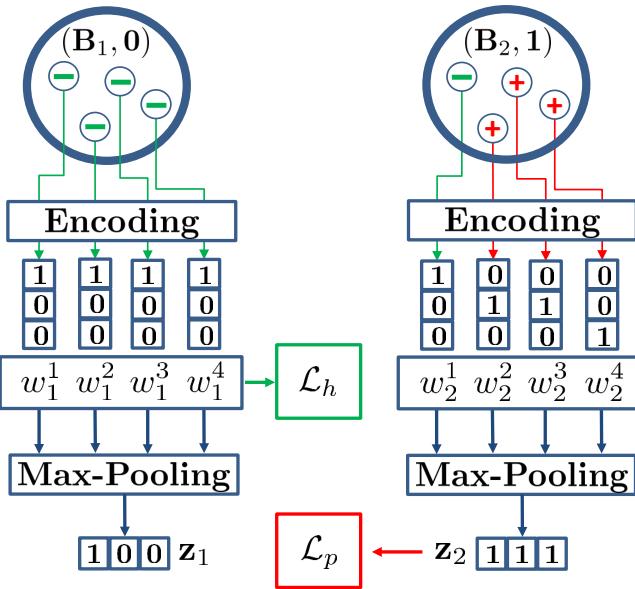


Fig. 4. A visual explanation of the interpretability-enhancing loss behavior. When the system receives a healthy image, every instance ideally contributes to the same codes in \mathcal{L}_h , while disease instances, only present on images showing signs of DR, appear as denser visual words contributing more to \mathcal{L}_p .

321 effect due to the contribution of all instances from the bag,
322 even when some of them may be irrelevant.

323 An alternative to avoid this effect is to perform a max-
324 pooling operation \mathbf{P} , as defined in eq. (2). In addition to
325 sharper mid-level representations, max-pooling matches better
326 the goal of DR detection. If no abnormal instance is found
327 on the mid-level representation associated to an image \mathbf{B} , it
328 should be declared as healthy. On the other hand, the
329 presence of a single microaneurysm or any other kind of
330 lesion is enough to classify the image as pathological. In our
331 case, max-pooling is implemented to accept the output codes
332 $\{\mathbf{w}^i, 1 \leq i \leq N(\mathbf{B})\}$ of the last hidden layer of \mathbf{U} and pool
333 the results into \mathbf{z} .

334 While training \mathbf{U} , the obtained mid-level representations
335 become the input of the second neural network \mathbf{D} , defined
336 in the same way as \mathbf{U} . The output of the final layer of \mathbf{D}
337 is supplied to a sigmoid activation unit, which produces a
338 single output containing the prediction of the system regarding
339 the presence or not of DR on the image \mathbf{B} from where the
340 instances were extracted. In training time, this prediction is
341 compared with the actual label of the image by means of a
342 cross-entropy loss penalizing inaccurate predictions:

$$\mathcal{L}_{class} = -\frac{1}{N} \sum_{i=1}^N l_i \log(\mathbf{D}(\mathbf{z}_i)) + (1 - l_i) \log(1 - \mathbf{D}(\mathbf{z}_i)), \quad (4)$$

343 where \mathbf{z}_i is the mid-level representation of the training image
344 \mathbf{B}_i , and l_i its corresponding ground-truth label. The weights
345 θ_U and θ_D of both networks are iteratively updated until
346 convergence by standard back-propagation with mini-batch
347 stochastic gradient descent, in order to minimize the error
348 given by eq. (4).

After jointly training \mathbf{U} and \mathbf{D} , given a new image \mathbf{B} , the
349 output of the proposed model is a prediction of the probability
350 p of \mathbf{B} being affected by DR, *i.e.*, $p = \mathbf{D}(\mathbf{P}(\mathbf{U}(\mathbf{B})))$.
351

B. A Strategy to Enforce Model Interpretability

Ideally, the mid-level representations \mathbf{z} obtained from pooling the encoded instances $\{\mathbf{w}^i, 1 \leq i \leq N(\mathbf{B})\}$ extracted from an image \mathbf{B} contain visually meaningful content. However, this behavior of the proposed model can be further enforced.
352
353
354
355
356
357

Since we know that healthy images only contain healthy instances, we can act at the instance level on codes \mathbf{w}^i extracted from healthy images. The goal in this case is to force the model to generate sparse mid-level representations. This can be accomplished by requiring that, when the label of an image \mathbf{B} is negative, \mathbf{U} uses few codes to encode all its instances. In this way, after pooling a set of codes from healthy instances, the resulting mid-level representation will necessarily be sparse.
358
359
360
361
362
363
364
365
366

In order to impose this behavior on the model, consider a healthy image \mathbf{B} , its set of instances $\{\mathbf{x}^i, 1 \leq i \leq N(\mathbf{B})\}$, and their corresponding codes $\{\mathbf{w}^i, 1 \leq i \leq N(\mathbf{B})\}$. We define the following quantity for healthy bags:
367
368
369
370

$$\mathcal{L}_h = \frac{1}{N} \sum_{i=1}^{N(\mathbf{B})} -\log(w_1^i), \quad (5)$$

which reaches a minimum whenever the code generated by \mathbf{U} is closer to the unitary vector $(1, 0, \dots, 0)$, since codes \mathbf{w}^i are normalized to sum up to 1.
371
372
373

On the other hand, if the image \mathbf{B} is pathological, it may contain both pathological and healthy instances. We cannot proceed in the same way and constrain the resulting mid-level representation from instance-level codes \mathbf{w}^i . However, we can still act at the bag level, by imposing that the mid-level representation \mathbf{z} is dense. In this case, we can define the following quantity for pathological bags:
374
375
376
377
378
379
380

$$\mathcal{L}_p = \frac{1}{M_L} \sum_{i=1}^{M_L} -\log(z_i), \quad (6)$$

which will be minimized whenever the mid-level representation \mathbf{z} is closer to the vector $(1, 1, \dots, 1)$. Since \mathbf{z} is the result of pooling codes coming from every instance in \mathbf{B} , this can only happen if the model encodes pathological instances with different visual words.
381
382
383
384
385
386
387
388

Finally, given a bag \mathbf{B} with corresponding label l , we define an interpretability-enhancement loss as the combination of both \mathcal{L}_h and \mathcal{L}_p :
389
390
391
392
393
394

$$\mathcal{L}_{int} = \alpha(1 - l)\mathcal{L}_h + \beta l \mathcal{L}_p. \quad (7)$$

Note that $l = 0$ whenever \mathbf{B} is healthy, whereas $l = 1$ for pathological images. In this way, each of the components of \mathcal{L}_{int} becomes active depending of bag-level information. Parameters α and β are positive real-valued hyper-parameters, weighting the contribution of each term. The global loss function that drives the learning of the entire system is simply
389
390
391
392
393
394

395 the addition of the binary classification loss defined in eq. (4)
 396 and the interpretability-enhancement loss of eq. (7):

$$\mathcal{L}_{global} = \mathcal{L}_{class} + \mathcal{L}_{int}. \quad (8)$$

397 A schematic representation illustrating the different situations
 398 the model may encounter, and the way in which the
 399 interpretability-enhancement loss in eq. (7) reacts to them, is
 400 shown in Fig. 4.

401 C. Implementation Details

402 In order to apply the proposed method to the DR detection
 403 problem, we need to decide on the feature extraction and
 404 description methods. In this case, Speeded Up Robust Features
 405 (SURF, [37]) were employed for both feature description and
 406 extraction as they have been shown to perform better than
 407 other description methods [29]. SURF is a scale and rotation
 408 invariant method that detects and describes interest points in an
 409 image. To better describe each interest point, 128 dimensional
 410 extended descriptors were computed. We used OpenCV's
 411 implementation of SURF [38] with default parameters and
 412 Theano [39] to implement the two neural networks U and
 413 D .

414 IV. EXPERIMENTAL EVALUATION

415 In this section, we provide experimental assessment of the
 416 performance of the proposed method when compared with
 417 other recent approaches. Performance is reported in terms
 418 of Area Under the ROC Curve (AUC) for the task of DR
 419 detection and DR referral, and we finally verify the enhanced
 420 interpretability of the proposed model, illustrating that it can
 421 effectively reveal the regions contributing to detect patholog-
 422 ical images.

423 A. DR Detection Performance Evaluation

424 We first evaluate the proposed DR detection technique on
 425 the publicly available Messidor dataset [40]. Messidor contains
 426 1200 color retinal fundus images acquired on three different
 427 French hospitals between 2005 and 2006. Images were ob-
 428 tained with TRC-NW6 non-mydriatic retinographs (Topcon,

TABLE II
 DR GRADING SCHEME FOR THE MESSIDOR DATASET
 MA = MICROANEURYSMS, HE = HARD EXUDATES, NV = NEO-VESSELS

DR	Description	Images
R0	$N_{MA} = 0$ and $N_{HE} = 0$	546
R1	$0 < N_{MA} \leq 5$ and $N_{HE} = 0$	153
R2	$5 < N_{MA} < 15$ and $0 < N_{HE} \leq 5$ and $N_{NV} = 0$	247
R3	$N_{MA} \geq 15$ or $N_{HE} \geq 5$ or $N_{NV} > 0$	254

Tokyo) with a 45° field of view, at a varying resolution of
 429 1440×960 , 2240×1488 and 2304×1536 . No image pre-
 430 processing was applied before extracting and describing the
 431 instances within these images.
 432

The clinical information associated to each image on Messidor
 433 consists of two labels indicating the grade of DR and the
 434 risk of macular edema, based on the presence and number of
 435 different types of lesions, see Table II. In order to build a DR
 436 presence label for each image, the provided values are merged
 437 in such a way that any image associated to a DR severity
 438 greater or equal than one is labeled as pathological, meaning
 439 that it contains early signs of DR, while only grade 0 images
 440 are considered as healthy. This resulted in a dataset on which
 441 546 images were labeled as normal and 654 as pathological.
 442

The Messidor dataset has been widely employed in the
 443 literature to assess the performance of DR detection and
 444 grading techniques. Some works approach the problem by
 445 designing a separate detector for each of possible lesion, and
 446 then applying it to Messidor images. The output of these lesion
 447 detectors is then combined with the set of rules in Table II in
 448 order to produce a DR presence/grade decision. However, it is
 449 worth noting that this approach requires the availability of an
 450 independent database containing pixel-wise ground-truth at the
 451 lesion level. This is precisely the challenge that our technique
 452 and other MIL-based methods try to overcome.
 453

A standard evaluation procedure was followed: 20% of the
 454 dataset was held-out for testing, while 65% was employed
 455 for training and 15% for validation. Hyper-parameters were
 456 found using random search [42], selecting the best values in
 457 terms of AUC on the validation set. Performance results of
 458

TABLE I
 PERFORMANCE COMPARISON OF DR DETECTION METHODS TESTED ON THE MESSIDOR DATASET.

Method	AUC	Observations
Red+Bright Lesion Detection [31]	88%	Shape, color, contrast features + kNN, combines [32] & [14]. Requires lesion annotations to be trained.
Ensemble [33]	88%	Ensembling of lesion detectors. Requires lesion annotations to be trained.
Red-Lesion Detection [13]	90%	Dynamic shape features + Random Forest classification. Requires lesion annotations to be trained.
DREAM [16]	90%	Ensembling of lesion detectors. Requires lesion annotations to be trained.
MIL Benchmark [34]	81%	Benchmarking of a set of 11 MIL techniques. Best result reported here.
Multiscale AM/FM [35]	84%	Frequency Analysis for Feature Extraction + Mahalanobis dist. for Classification
AM/FM - SVM [36]	86%	Frequency Analysis for Feature Extraction + SVM for Classification
MIL for DR detection [28]	88%	BoVW with separate Encoding+Classification training.
Data-Mined Context [30]	89%	GFTT+SIFT features. Mines contextual data from patient's record, including text.
Ours	90%	Proposed Approach: MIL with joint Encoding+Classification training.

TABLE IV
PERFORMANCE COMPARISON OF DR DETECTION ON THE DR1 DATASET.

Method	AUC	Observations
Multi-Lesion fusion [43]	84%	Lesion detection + meta-SVM.
Ours	93%	Proposed Approach.

459 the proposed technique are shown in Table I in terms of AUC,
 460 together with the performance obtained in the same dataset by
 461 different state-of-the-art techniques. We include both methods
 462 trained on independent datasets for the task of lesion detection
 463 and methods that learn directly from the image-level ground-
 464 truth in order to predict DR detection.

465 The results on Table I lead to several conclusions. First, the
 466 proposed method achieves a superior performance **comparing**
 467 **to the other** techniques that have been trained without access
 468 to pixel-level lesion annotations. It is particularly interesting
 469 to note that other MIL-based techniques such as [28], or the
 470 best of the DR detection techniques reported in [34], obtain a
 471 lower AUC. The main difference between all these methods
 472 and the technique introduced in this paper is the propagation of
 473 the bag-level labels until the encoding process, which directly
 474 benefits from this information in order to produce more
 475 useful mid-level representations, leading to a better detection
 476 performance. Second, performance of methods trained with
 477 lesion-level ground-truth is comparable but not superior to
 478 the introduced technique. This means that the proposed MIL-
 479 based approach for DR detection can effectively make use of
 480 local information on the image to the same extent as these
 481 techniques, but without having explicit access to it.

482 It should be noted that other recent approaches based on
 483 Deep Convolutional Neural Networks (CNN) have been tested
 484 with great success on the Messidor dataset [6], achieving
 485 even larger AUC values without the need of lesion-level
 486 information. However, this study proposes a model trained on
 487 an external large dataset of retinal images. This private dataset
 488 contained 118 175 training images, independently labeled 3
 489 to 7 times by a panel of 54 ophthalmologists. Moreover,
 490 the output of this kind of CNN-based models typically lacks
 491 interpretability, which may hinder the predisposition of doctors
 492 towards its acceptance in a regular clinical workflow. The
 493 method introduced in this paper addresses both issues by
 494 leveraging as much information as possible from a moderate-
 495 size dataset, while enforcing the interpretable behavior of the
 496 model, as illustrated in section IV-C.

497 In order to test if the proposed method generalized to
 498 different datasets, we also tried to detect DR as the presence of

any single lesion in the DR1 dataset, introduced in [41]. This
 499 dataset contains 1077 retinal images captured with a TRC-50X
 500 (Topcon Inc., Tokyo, Japan) mydriatic camera with a 45° field
 501 of view and an average resolution of 640 × 1077 pixels. From
 502 all the images, 595 were classified as containing no sign of
 503 DR and 482 as showing pathological signs. In this case, we are
 504 only aware of a work addressing the task of DR detection on
 505 this dataset [43]. Since DR1 contains ground-truth regarding
 506 the presence of different lesions within each pathological
 507 image, DR detection is achieved by training separate detectors
 508 for each of them and then fusing the results with a meta-
 509 classifier. Performance comparison **with the results in [43] is**
 510 **shown** in Table IV. In this case, the proposed technique clearly
 511 outperforms the method proposed in [43], which confirms the
 512 generality of our approach.

B. DR Referral Performance Evaluation

513 It has been argued in [41] that the presence of a given lesion
 514 may not be enough to make a decision on the need to refer a
 515 patient for further examination. Table II contains a set of rules
 516 designed in order to make such a decision, but it may not cover
 517 all the signs an expert ophthalmologist takes into account when
 518 recommending further examination of a patient. For instance,
 519 in [44], referral is defined as having a DR grade above mild
 520 non-proliferative (R1) and/or macular edema. Lesion location
 521 may also impact the clinical decision regarding referral.

522 In order to assess the performance of the proposed method
 523 in terms of DR referral prediction, we employ the publicly
 524 available DR2 dataset. This dataset was introduced in [41], and
 525 it is composed of 520 retinal fundus images, from which 337
 526 images were categorized by two independent ophthalmologists
 527 as not requiring referral, and 98 were deemed to require
 528 referral within one year by a specialist. It is important to note
 529 that while labeling the images, the experts were required to
 530 categorize them ignoring specific lesions and considering only
 531 if the image should lead to referral. The medical specialists
 532 based their decision on any reason considered to be clin-
 533 ically relevant, not only on the presence of particular lesions.
 534 DR2 images were acquired with a TRC-NW8 (Topcon Inc.,
 535 Tokyo, Japan) nonmydriatic retinal camera, and they all have
 536 867 × 575 pixel resolution and a 45° field of view.

537 Several methods have been proposed in the past for DR
 538 referral prediction [36], [44], [45]. Unfortunately, there is not
 539 a standard definition of referral, which results in different
 540 problems of varying difficulty being solved. In [45] images
 541 containing signs of macular edema were considered as re-
 542 ferable, while in [36] the adopted definition was the presence
 543

TABLE III
PERFORMANCE COMPARISON OF DR DETECTION METHODS TESTED ON THE DR2 DATASET.

Method	AUC	Observations
Pires 2013 [41]	93	Separate Lesion Detectors + Meta-Classifier. Requires weak lesion information to be trained.
Pires 2014 [29]	94%	Separate Lesion Detectors + Meta-Classifier. Requires weak lesion information to be trained.
Pires 2017 [26]	96%	Bypasses Lesion Detection. BossaNova and Fisher Vector features + BoVW.
Ours	96%	Proposed Approach.

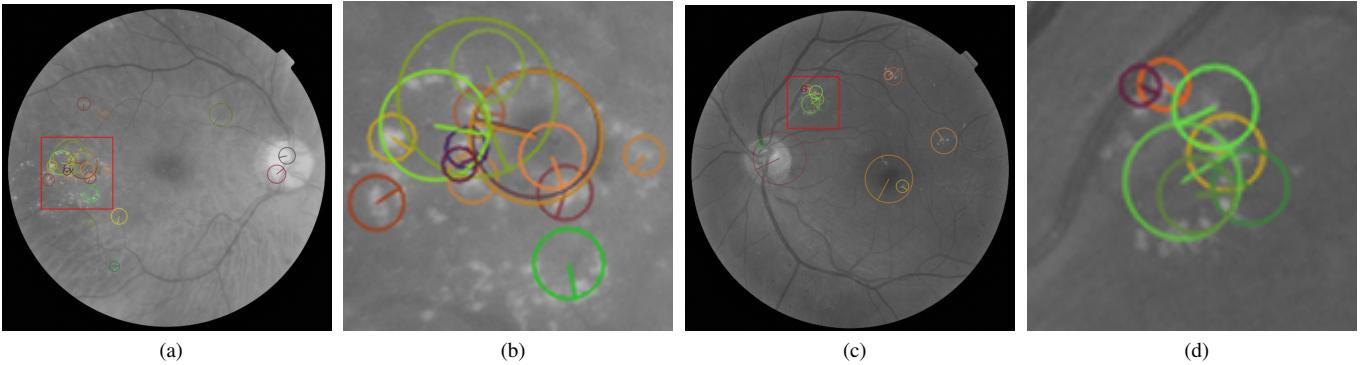


Fig. 5. Instances on two pathological images that contributed to the decision produced by the proposed system. In this case, it can be appreciated that most of the SURF keypoints are located on top of bright lesions. Best viewed in color.

of signs of high DR grade or high risk of macular edema. In order to be able to properly compare the proposed approach in a fair manner, we select those **studies** that were tested on DR2, since it contains direct referral opinion from medical experts. In this case, [41] proposed a solution consisting of training individual lesion detectors, and employing the resulting decision scores in order to train a meta-classifier to predict referral. The individual detectors consist on a variant of BoVW with more advanced pooling and encoding operations. In a later work [29], this scheme was improved by means of a semi-soft encoding strategy, with results outperforming those of [41].

It is worth noting that these methods, even if being MIL-based approaches, still need to be trained with weak lesion-level information. In this case, information regarding which kind of lesion is present in an image, but not the exact location and **its delineation, is used**. In order to train these techniques, the DR1 dataset described in the previous section was used. In contrast, the technique proposed in this paper was trained only on DR2, with no other information than the need for referral of each image. This characteristic is shared by another recent technique introduced in [26]. In that work, it was effectively shown that DR referral could be predicted without the need for explicit lesion detection. However, the proposed method still presents a separate visual dictionary construction and classifier

training.

Performance results for DR referable predictions in terms of AUC is presented in table III. It can be observed that the proposed technique improves or matches the performance of previously reported methods also in the task of DR referral. The arguments suggested in [26] about the possibility of training a referral prediction system without the need of explicitly building separate lesion detectors **are confirmed** by these results. We can conclude that both the technique proposed in this work and the one introduced in [26] obtain superior performance than lesion-detection based techniques, confirming the validity of this approach. **It is important to notice, however, that the method from [26] only addresses DR referral, while the technique presented here is tested both in DR detection and referral. Moreover, the results obtained by our technique are better interpretable than those produced by [26].** In the next section we analyze this aspect of the proposed model.

C. Interpretability of the Model

One of the most relevant features of the DR detection system proposed in this paper is its enhanced interpretability, allowed by the joint minimization of the two loss functions in eqs. (4) and (7). This enables us to explain which where the instances within the images that most likely caused the model

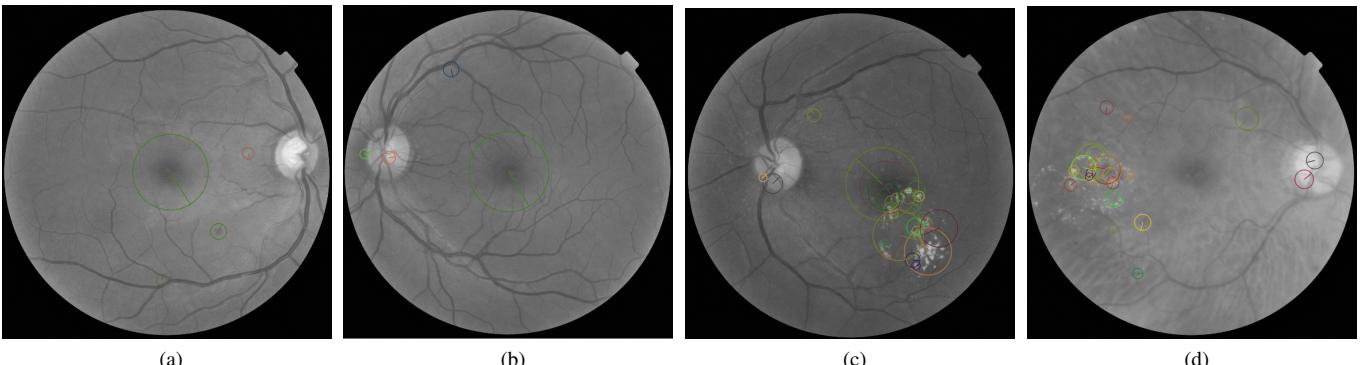


Fig. 6. SURF keypoints associated to image instances considered by the proposed model in order to produce a decision on DR presence. (a) and (b) depict healthy images, on which fewer instances were taken into account. On the contrary, (c) and (d) show pathological images, on which a greater amount of instances are considered to reach a decision. Best viewed in color.

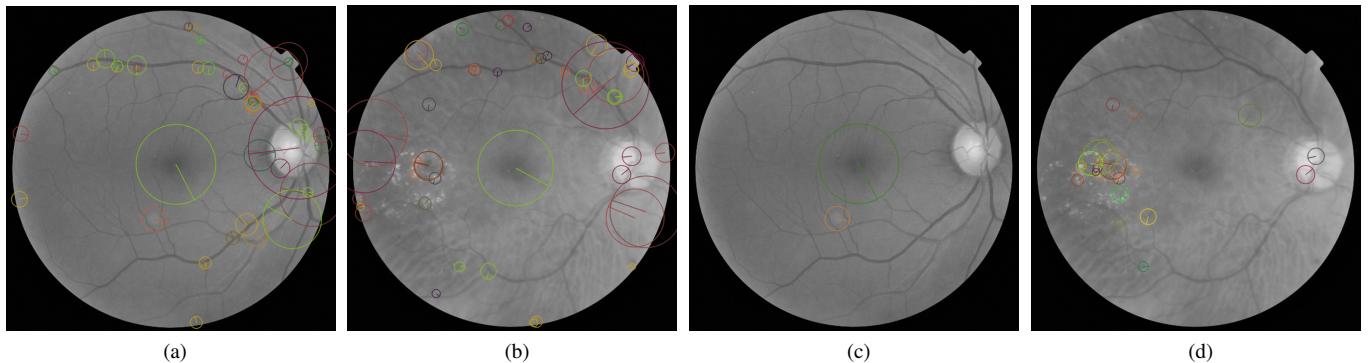


Fig. 7. Comparison between results produced by a model trained without the interpretability-enhancement loss - (a), (b) - and after adding it - (c), (d). In this case, (a) and (c) show a healthy image, while (b) and (d) show a pathological example. Best viewed in color.

594 to reach the produced decision. To experimentally demonstrate
 595 this aspect of the model, **a first example of this behavior on**
 596 **pathological images from the DR2 dataset is shown in Fig. 5.**
 597 In Figs. (5a) and (5c), SURF keypoints contributing to the
 598 resulting mid-level representations of these images are de-
 599 picted. We can clearly see how most of the selected instances
 600 correspond to keypoints extracted from bright lesions, while
 601 only few keypoints are related to instances that are typically
 602 present on both normal and pathological retinal images, such
 603 as the macula or the optic disc. Zoomed-in details are also
 604 shown in Figs. (5b) and (5d) to better present this observation.

605 **A second experiment was run to better illustrate that the**
 606 **model trained with the interpretability-enhancement loss be-**
 607 **haves as desired.** The loss function introduced in eq. (7) aims
 608 at promoting a sparse mid-level representation for healthy
 609 images, on which few instances are ideally considered, while
 610 in the case of pathological images, the produced mid-level
 611 representations are expected to be denser. This should translate
 612 into more SURF keypoints appearing when considering patho-
 613 logical examples. This is visually verified in Fig. 6. There it
 614 can be observed that **fewer** keypoints were taken into account
 615 when reaching a decision regarding a healthy image, see Figs.
 616 (6a) and (6b), than when a pathological image was considered,
 617 as shown in Figs. (6c) and (6d).

618 To further verify that the loss term in eq. (7) effectively
 619 contributes to the explainability of the model's decision, we
 620 trained a separate classifier by minimizing only the loss
 621 function of eq. (4), without including the interpretability-
 622 enhancement loss term. Both results are visually compared
 623 in Fig. 7. It can be seen how the extra loss term leads to
 624 more interpretable results by a better identification of the
 625 pathological instances. Only a fraction of the input instances
 626 are used by the model to produce a decision, while irrelevant
 627 instances are filtered out. It is worth noting that when the
 628 interpretability-enhancement loss term was not included, the
 629 model considered roughly the same number of keypoints on
 630 normal and pathological images, as shown in Figs. (7a) and
 631 (7b). However, when the global loss in eq. (7) is minimized,
 632 the resulting model considers substantially less keypoints in a
 633 healthy example than in a pathological image in order to make
 634 a decision, as can be **observed** in Figs. (7c) and (7d).

V. CONCLUSIONS AND FUTURE WORK

In this paper, a new Diabetic Retinopathy (DR) detection system has been presented, based on the Multiple-Instance Learning framework. The method can learn from weak information regarding only the presence or absence of disease to formulate predictions on new images based on implicit local information. The main novelty of the proposed model with respect to previously existing MIL-based DR detection systems is a joint-learning scheme in which the encoding and the classification stages are connected. Thanks to this approach, the mid-level representations generated by the model are optimized to improve DR detection accuracy. Furthermore, a novel strategy to enforce the interpretability of the resulting predictions has been introduced, resulting in a better understanding of the output of the model. Performance comparisons against other recent DR detection and DR referral techniques give advantage to the proposed technique, confirming previous observations stating that weak expert labels (at the image level only) can be leveraged to produce accurate predictions without the need of pixel-level information related to the different lesions indicating the presence of DR.

The developed technique achieves good performance, but further improvements can be achieved. Speeded-Up Robust Features (SURF) were employed in this work to locate and describe instances within retinal images. Even if the proposed technique jointly optimizes the encoding and the classification stages of the model, instance location and description may be included in the same global optimization process. This can be achieved with an end-to-end system in which the most appropriate image representation for the task of DR detection is also learned by using a Deep Convolutional Neural Network. Further work will involve exploring this direction of research, in order to obtain higher performance in terms of DR detection, as well as extending the approach to predicting different levels of DR severity.

ACKNOWLEDGMENT

This work is financed by the ERDF – European Regional Development Fund through the Operational Programme for Competitiveness and Internationalisation - COMPETE 2020 Programme, and by National Funds through the FCT -

675 Fundação para a Ciência e a Tecnologia within project CMUP-
 676 ERI/TIC/0028/2014.

REFERENCES

- [1] R. Gargeya and T. Leng, "Automated Identification of Diabetic Retinopathy Using Deep Learning," *Ophthalmology*, vol. 124, no. 7, pp. 962–969, Jul. 2017.
- [2] D. M. Squirrell and J. F. Talbot, "Screening for diabetic retinopathy," *Journal of the Royal Society of Medicine*, vol. 96, no. 6, pp. 273–276, Jun. 2003.
- [3] P. H. Scanlon, "The English National Screening Programme for diabetic retinopathy 2003–2016," *Acta Diabetologica*, vol. 54, no. 6, pp. 515–525, Jun. 2017.
- [4] L. P. Daskivich, C. Vasquez, C. Martinez, C.-H. Tseng, and C. M. Mangione, "Implementation and Evaluation of a Large-Scale Teleretinal Diabetic Retinopathy Screening Program in the Los Angeles County Department of Health Services," *JAMA Internal Medicine*, vol. 177, no. 5, pp. 642–649, May 2017.
- [5] J. Beagley, L. Guariguata, C. Weil, and A. A. Motala, "Global estimates of undiagnosed diabetes in adults," *Diabetes Research and Clinical Practice*, vol. 103, no. 2, pp. 150–160, Feb. 2014.
- [6] V. Gulshan, L. Peng, M. Coram, M. C. Stumpe, D. Wu, A. Narayanaswamy, S. Venugopalan, K. Widner, T. Madams, J. Cuadros, R. Kim, R. Raman, P. C. Nelson, J. L. Mega, and D. R. Webster, "Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs," *JAMA*, vol. 316, no. 22, pp. 2402–2410, Dec. 2016.
- [7] A. D. Fleming, S. Philip, K. A. Goatman, G. J. Prescott, P. F. Sharp, and J. A. Olson, "The evidence for automated grading in diabetic retinopathy screening," *Current Diabetes Reviews*, vol. 7, no. 4, pp. 246–252, Jul. 2011.
- [8] E. Soto-Pedre, A. Navea, S. Millan, M. C. Hernaez-Ortega, J. Morales, M. C. Desco, and P. Pérez, "Evaluation of automated image analysis software for the detection of diabetic retinopathy to reduce the ophthalmologists' workload," *Acta Ophthalmologica*, vol. 93, no. 1, pp. e52–56, Feb. 2015.
- [9] P. Costa and A. Campilho, "Convolutional bag of words for diabetic retinopathy detection from eye fundus images," in *2017 Fifteenth IAPR International Conference on Machine Vision Applications (MVA)*, May 2017, pp. 165–168.
- [10] M. Niemeijer, B. van Ginneken, M. J. Cree, A. Mizutani, G. Quellec, C. I. Sanchez, B. Zhang, R. Hornero, M. Lamard, C. Muramatsu, X. Wu, G. Cazuguel, J. You, A. Mayo, Q. Li, Y. Hatanaka, B. Cochener, C. Roux, F. Karray, M. Garcia, H. Fujita, and M. D. Abramoff, "Retinopathy online challenge: automatic detection of microaneurysms in digital color fundus photographs," *IEEE Transactions on Medical Imaging*, vol. 29, no. 1, pp. 185–195, Jan. 2010.
- [11] C. I. Sánchez, M. Niemeijer, M. S. A. S. Schulten, M. Abràmoff, and B. v. Ginneken, "Improving hard exudate detection in retinal images through a combination of local and contextual information," in *2010 IEEE International Symposium on Biomedical Imaging: From Nano to Macro*, Apr. 2010, pp. 5–8.
- [12] R. Srivastava, D. W. K. Wong, L. Duan, J. Liu, and T. Y. Wong, "Red lesion detection in retinal fundus images using Frangi-based filters," in *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, Aug. 2015, pp. 5663–5666.
- [13] L. Seoud, T. Hurtut, J. Chelbi, F. Cheriet, and J. M. P. Langlois, "Red Lesion Detection Using Dynamic Shape Features for Diabetic Retinopathy Screening," *IEEE Transactions on Medical Imaging*, vol. 35, no. 4, pp. 1116–1126, Apr. 2016.
- [14] M. Niemeijer, B. van Ginneken, S. R. Russell, M. S. A. Suttorp-Schulten, and M. D. Abràmoff, "Automated detection and differentiation of drusen, exudates, and cotton-wool spots in digital color fundus photographs for diabetic retinopathy diagnosis," *Investigative Ophthalmology & Visual Science*, vol. 48, no. 5, pp. 2260–2267, May 2007.
- [15] Ujjwal, K. S. Deepak, A. Chakravarty, and J. Sivaswamy, "Visual saliency based bright lesion detection and discrimination in retinal images," in *2013 IEEE 10th International Symposium on Biomedical Imaging*, Apr. 2013, pp. 1436–1439.
- [16] S. Roychowdhury, D. D. Koozekanani, and K. K. Parhi, "DREAM: Diabetic Retinopathy Analysis Using Machine Learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 18, no. 5, pp. 1717–1728, Sep. 2014.
- [17] I. N. Figueiredo, S. Kumar, C. M. Oliveira, J. D. Ramos, and B. Enquist, "Automated lesion detectors in retinal fundus images," *Computers in Biology and Medicine*, vol. 66, pp. 47–65, Nov. 2015.
- [18] S. Andrews, I. Tsochantarisidis, and T. Hofmann, "Support Vector Machines for Multiple-instance Learning," in *Proceedings of the 15th International Conference on Neural Information Processing Systems*, ser. NIPS'02. Cambridge, MA, USA: MIT Press, 2002, pp. 577–584.
- [19] P. Viola, J. C. Platt, and C. Zhang, "Multiple Instance Boosting for Object Detection," in *Proceedings of the 18th International Conference on Neural Information Processing Systems*, ser. NIPS'05. Cambridge, MA, USA: MIT Press, 2005, pp. 1417–1424.
- [20] Y. Chen, J. Bi, and J. Z. Wang, "MILES: Multiple-Instance Learning via Embedded Instance Selection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 12, pp. 1931–1947, Dec. 2006.
- [21] J. Sivic and A. Zisserman, "Video Google: a text retrieval approach to object matching in videos," in *Proceedings Ninth IEEE International Conference on Computer Vision*, Oct. 2003, pp. 1470–1477 vol.2.
- [22] G. Quellec, G. Cazuguel, B. Cochener, and M. Lamard, "Multiple-Instance Learning for Medical Image and Video Analysis," *IEEE Reviews in Biomedical Engineering*, vol. PP, no. 99, pp. 1–1, 2017.
- [23] V. Cheplygina, L. Sørensen, D. M. J. Tax, J. H. Pedersen, M. Loog, and M. d. Bruijne, "Classification of COPD with Multiple Instance Learning," in *2014 22nd International Conference on Pattern Recognition*, Aug. 2014, pp. 1508–1513.
- [24] Y. Xu, J. Y. Zhu, E. Chang, and Z. Tu, "Multiple clustered instance learning for histopathology cancer image classification, segmentation and clustering," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2012, pp. 964–971.
- [25] T. Tong, R. Wolz, Q. Gao, R. Guerrero, J. V. Hajnal, and D. Rueckert, "Multiple instance learning for classification of dementia in brain MRI," *Medical Image Analysis*, vol. 18, no. 5, pp. 808–818, Jul. 2014.
- [26] R. Pires, S. Avila, H. F. Jelinek, J. Wainer, E. Valle, and A. Rocha, "Beyond Lesion-Based Diabetic Retinopathy: A Direct Approach for Referral," *IEEE Journal of Biomedical and Health Informatics*, vol. 21, no. 1, pp. 193–200, Jan. 2017.
- [27] S. Manivannan, C. Cobb, S. Burgess, and E. Trucco, "Subcategory Classifiers for Multiple-Instance Learning and Its Application to Retinal Nerve Fiber Layer Visibility Classification," *IEEE Transactions on Medical Imaging*, vol. 36, no. 5, pp. 1140–1150, May 2017.
- [28] G. Quellec, M. Lamard, M. D. Abràmoff, E. Decencière, B. Lay, A. Erginay, B. Cochener, and G. Cazuguel, "A multiple-instance learning framework for diabetic retinopathy screening," *Medical Image Analysis*, vol. 16, no. 6, pp. 1228–1240, Aug. 2012.
- [29] R. Pires, H. F. Jelinek, J. Wainer, E. Valle, and A. Rocha, "Advancing Bag-of-Visual-Words Representations for Lesion Classification in Retinal Images," *PLOS ONE*, vol. 9, no. 6, p. e96814, Jun. 2014.
- [30] G. Quellec, M. Lamard, A. Erginay, A. Chabouis, P. Massin, B. Cochener, and G. Cazuguel, "Automatic detection of referral patients due to retinal pathologies through data mining," *Medical Image Analysis*, vol. 29, pp. 47–64, Apr. 2016.
- [31] C. I. Sánchez, M. Niemeijer, A. V. Dumitrescu, M. S. A. Suttorp-Schulten, M. D. Abràmoff, and B. v. Ginneken, "Evaluation of a Computer-Aided Diagnosis System for Diabetic Retinopathy Screening on Public Data," *Investigative Ophthalmology & Visual Science*, vol. 52, no. 7, pp. 4866–4871, Jun. 2011.
- [32] M. Niemeijer, B. v. Ginneken, J. Staal, M. S. A. Suttorp-Schulten, and M. D. Abràmoff, "Automatic detection of red lesions in digital color fundus photographs," *IEEE Transactions on Medical Imaging*, vol. 24, no. 5, pp. 584–592, May 2005.
- [33] B. Antal and A. Hajdu, "An Ensemble-Based System for Microaneurysm Detection and Diabetic Retinopathy Grading," *IEEE Transactions on Biomedical Engineering*, vol. 59, no. 6, pp. 1720–1726, Jun. 2012.
- [34] M. Kandemir and F. A. Hamprecht, "Computer-aided diagnosis from weak supervision: a benchmarking study," *Computerized Medical Imaging and Graphics: The Official Journal of the Computerized Medical Imaging Society*, vol. 42, pp. 44–50, Jun. 2015.
- [35] C. Agurto, V. Murray, E. Barriga, S. Murillo, M. Pattichis, H. Davis, S. Russell, M. Abràmoff, and P. Soliz, "Multiscale AM-FM Methods for Diabetic Retinopathy Lesion Detection," *IEEE Transactions on Medical Imaging*, vol. 29, no. 2, pp. 502–512, Feb. 2010.
- [36] E. S. Barriga, V. Murray, C. Agurto, M. Pattichis, W. Bauman, G. Zamora, and P. Soliz, "Automatic system for diabetic retinopathy screening based on AM-FM, partial least squares, and support vector machines," in *2010 IEEE International Symposium on Biomedical Imaging: From Nano to Macro*, Apr. 2010, pp. 1349–1352.
- [37] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-Up Robust Features (SURF)," *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, Jun. 2008.
- [38] "OPENCV, Open Source Computer Vision Library," 2015. [Online]. Available: <https://github.com/itseez/opencv>

- 827 [39] Theano Development Team, "Theano: A Python framework for
828 fast computation of mathematical expressions," *arXiv e-prints*, vol.
829 abs/1605.02688, May 2016.
- 830 [40] E. Decencière, X. Zhang, G. Cazuguel, B. Lay, B. Cochener, C. Trone,
831 P. Gain, R. Ordonez, P. Massin, A. Erginay, B. Charton, and J.-C.
832 Klein, "Feedback on a publicly distributed image database: the Messidor
833 database," *Image Analysis & Stereology*, vol. 33, no. 3, pp. 231–234,
834 Aug. 2014.
- 835 [41] R. Pires, H. F. Jelinek, J. Wainer, S. Goldenstein, E. Valle, and A. Rocha,
836 "Assessing the Need for Referral in Automatic Diabetic Retinopathy
837 Detection," *IEEE transactions on bio-medical engineering*, vol. 60,
838 no. 12, pp. 3391–3398, Dec. 2013.
- 839 [42] J. Bergstra and Y. Bengio, "Random Search for Hyper-Parameter Opti-
840 mization," *Journal of Machine Learning Research*, vol. 13, no. Feb, pp.
841 281–305, 2012.
- 842 [43] H. F. Jelinek, R. Pires, R. Padilha, S. Goldenstein, J. Wainer, T. Bosso-
843 maier, and A. Rocha, "Data fusion for multi-lesion Diabetic Retinopathy
844 detection," in *2012 25th IEEE International Symposium on Computer-
845 Based Medical Systems (CBMS)*, Jun. 2012, pp. 1–4.
- 846 [44] M. D. Abràmoff, J. C. Folk, D. P. Han, J. D. Walker, D. F. Williams,
847 S. R. Russell, P. Massin, B. Cochener, P. Gain, L. Tang, M. Lamard,
848 D. C. Moga, G. Quellec, and M. Niemeijer, "Automated Analysis of
849 Retinal Images for Detection of Referable Diabetic Retinopathy," *JAMA
850 Ophthalmology*, vol. 131, no. 3, pp. 351–357, Mar. 2013.
- 851 [45] K. S. Deepak and J. Sivaswamy, "Automatic assessment of macular
852 edema from color retinal images," *IEEE Transactions on Medical
853 Imaging*, vol. 31, no. 3, pp. 766–776, Mar. 2012.