

# Real-time Lexicon-Free Scene Text Retrieval

Andrés Mafla, Rubèn Tito, Lluís Gómez, Marçal Rusiñol, Ernest Valveny,  
Dimosthenis Karatzas

*Computer Vision Center, Universitat Autònoma de Barcelona. Edifici O, Campus UAB,  
08193 Bellaterra (Cerdanyola) Barcelona, Spain.*

*E-mail: {amafla,rperez,lgomez,marcal,ernest,dimos}@cvc.uab.cat*

---

## Abstract

In this work, we address the task of scene text retrieval: given a text query, the system must return all images containing the queried text. The proposed model uses a single shot CNN architecture that predicts bounding boxes and builds a compact representation of spotted words. In this way, this problem can be modeled as a nearest neighbor search of the textual representation of a query over the outputs of the CNN collected from the totality of an image database. Our experiments demonstrate that the proposed model outperforms previous state-of-the-art, while offering a significant increase in processing speed and unmatched expressiveness with samples never seen at training time. Several experiments to assess the robustness of the model are conducted as well as an application of real-time text spotting in videos.

*Keywords:* Image retrieval, Scene text detection, Scene text recognition, Word spotting, Convolutional Neural Networks, Region Proposals Networks, PHOC.

---

# Real-time Lexicon-Free Scene Text Retrieval

Andrés Mafla, Rubèn Tito, Lluís Gómez, Marçal Rusiñol, Ernest Valveny,  
Dimosthenis Karatzas

*Computer Vision Center, Universitat Autònoma de Barcelona. Edifici O, Campus UAB,  
08193 Bellaterra (Cerdanyola) Barcelona, Spain.*

*E-mail: {amafla,rperez,lgomez,marcal,ernest,dimos}@cvc.uab.cat*

---

---

## 1. Introduction

The development of language is one of the most influential inventions of humankind that allows the communication of abstract and complex ideas. Similarly, written text permits this set of complex ideas to be depicted in an explicit and semantic manner. As it is shown by several authors [1, 2, 3], there is a big percentage of media that contains text, especially in urban scenarios and documents. Adding this to the fact that there is ample availability of data and the importance of text, it becomes essential to develop and refine algorithms that exploit the richness of textual information found in images and video. Leveraging text in scene imagery allows the emergence of tasks such as image retrieval [4, 5], scene understanding [6, 7], instant translation [8, 9], human-computer interaction, robot navigation [10, 11], assisted reading for the visually-impaired [12, 13] and industrial automation [14, 15]. In the previous years significant advances have been accomplished, particularly since the introduction of AlexNet [16], architecture that won the ILSVRC2012 [17] contest by using deep learning techniques. Text spotting has been diverging from older approaches that used hand-crafted features

24 towards current ones that employ automatic feature learning by exploiting  
25 deep learning methodologies [12, 18]. Nonetheless, text spotting is not a triv-  
26 ial task and remains as an open problem in the research community. Putting  
27 aside the complexity of spotting text in the wild, the importance that text  
28 encompasses is given by the high level semantic and explicit information,  
29 which can not be leveraged by using visual cues alone. For example, there is  
30 a high degree of complexity involved in labelling images without considering  
31 the text found in them, even for humans. This effect is evident in Figure 1,  
32 in which the storefronts alone can belong to a wide plethora of businesses,  
33 but the exact label can be inferred if and only if the text contained is read  
34 and leveraged appropriately. Research conducted by Movshovitz *et al.* [19]  
35 showed that while training a shop classifier, the proposed model ended up  
36 learning and interpreting textual information as the only way of differentiat-  
37 ing between diverse businesses. The described effect is evident and addressed  
38 explicitly in later works conducted by [20, 21], which focuses on fine-grained  
39 classification of storefronts and bottles respectively. Additional tasks that  
40 require integration of textual and visual information to generate a common  
41 domain knowledge have been proposed such as in [6, 7], which opens up new  
42 research paths.

43 Closely related to our work, Mishra *et al.* [22] proposed the task of scene  
44 text retrieval. The input to the system is a text query, which the system  
45 must employ to return all the images that contain the queried text. This  
46 task requires systems that are robust enough to perform fast word spotting  
47 while at the same time holding the capacity of generalizing out of dictionary  
48 queries never seen before. An intuitive approach to tackle such a problem

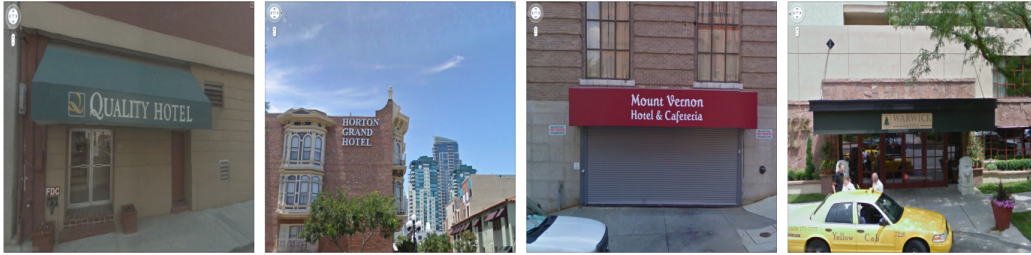


Figure 1: The visual appearance of different business places in images can be extremely variable. It seems impossible to correctly label them without reading the text within them. Our scene text retrieval method returns all the images shown here within the top-10 ranked results among more than 10,000 distractors for the text query “hotel”.

49 is to make use of state of the art reading systems, and use the output pre-  
50 dictions of it to find the closest match with the given query. However, as it  
51 is shown by [22], such attempts commonly have low performance caused by  
52 limitations in end to end reading systems. On one hand, end to end reading  
53 systems are evaluated on recognition, a different task that focuses on achiev-  
54 ing high precision scores, often using a specific language dictionary [23] or as  
55 it is proposed by [24, 13] a short dictionary per image. On the other hand, a  
56 retrieval system requires a large number of proposals (high recall) which can  
57 be beneficial at the moment of finding close matching detections when com-  
58 pared to a query. It is worth noting that end to end reading systems usually  
59 consist of at least two clearly defined stages that employ the encoder-decoder  
60 paradigm. The pipeline comprised by these two stages, more often that not  
61 are slow at the moment of generating predictions of the text contained in  
62 an image. This existing time constraint hinders the use of such algorithms  
63 in real-time scenarios or at the moment of indexing large scale collections of  
64 images and documents.

65 In order to exploit the particular requirements that need to be addressed



66 by a retrieval system, we propose in this work a real-time, high-performance  
67 word spotting method that detects and recognizes text in a single calculation  
68 of a Fully Convolutional Neural Network (FCNN). The proposed architec-  
69 ture is based on the YOLO model [25, 26], a widely used single shot object  
70 detector which in our case is employed to construct a PHOC (Pyramidal His-  
71 togram Of Characters) [27, 28] predictor. By employing this methodology,  
72 our model is able to perform text detection and recognition in a single cal-  
73 culation thus making it suitable for real time applications or to index large  
74 scale image collections at an unmatched speed.

75 The main contributions of using the proposed model, as it is shown in our  
76 previous work [29] are: firstly, the usage of a layout comprised by an end-  
77 to-end jointly trainable FCNN. Secondly, the usage of the PHOC as a word  
78 representation instead of a direct word classification over a closed dictionary.  
79 Thus, providing an elegant mechanism to generalize to any text string, al-  
80 lowing the method to tackle efficiently out-of-dictionary queries. Lastly, due  
81 to its design, the adoption of this method achieves unmatched speed when  
82 processing images to construct a compact representations of the recognized  
83 text instances. As an extension to the preceding research, in this work we  
84 analyze deeply the capacity of dealing with out-of-vocabulary queries of our  
85 model by conducting exhaustive experiments performed in two multi-lingual  
86 datasets. These experiments prove that the proposed method is able to suc-  
87 cessfully apply knowledge transfer acquired at training time to construct  
88 word representations of previously unseen text samples at inference time. As  
89 an additional section we present supplementary experiments and provide an  
90 analysis of the system under different kinds of imperfect image conditions

91 such as rotation, blur, occlusion and compression, experiments that confirm  
92 the robustness of the proposed architecture. Lastly, we propose an applica-  
93 tion of real-time text spotting on video, in which the model needs to confirm  
94 its robustness to noise and distortions while at the same time maintaining  
95 its characteristic high processing speed.

## 96 **2. Related Work**

97 In the past years, several advances in Deep Learning have been accom-  
98 plished due to data availability and computing power [3], allowing deep learn-  
99 ing models to surpass several benchmarks in a wide range of tasks. The main  
100 advantage of using deep learning methodologies is the possibility of automatic  
101 feature learning, rather than hand-crafted ones. Most literature [18, 12] di-  
102 vide the existing methods as: text detection, text recognition, end-to-end  
103 systems. Other applications such as fine-grained classification, image under-  
104 standing and image retrieval are briefly described in the upcoming sections.

### 105 *2.1. Scene Text Detection*

106 Initial deep learning methodologies employed several steps to produce  
107 proposals. In the work presented by [30], a CNN is used to predict if a  
108 given pixel belongs to a character, forms part of a text region and its ori-  
109 entation. Yao *et al.* [31] propose a CNN that outputs text proposals, which  
110 are filtered by separating different text instances by employing a semantic  
111 segmentation model. Later works focus on simplifying the pipeline and thus  
112 improving speed and training of models. These models usually follow a two  
113 step pipeline that comprise of an end-to-end trainable detection network and  
114 a post processing step. The work presented by [32, 33] named Textboxes,

115 adopts a modified version of a popular object recognition model named Single  
116 Shot Detector [34]. It employs modified anchor boxes to regress the ground  
117 truth boxes followed by a non-maximum suppression step (NMS). A per-  
118 formance focused approach is given by EAST [35], which upsamples feature  
119 maps gradually and uses [36] as the network backbone, and outputs a per  
120 pixel word or text line prediction followed by a NMS step.

121 Inspired by the object detection framework proposed by R-CNN [37, 38, 39],  
122 ample research has been conducted. The common approach consists of a Re-  
123 gion Proposal Network (RPN) that produces candidate text regions, which  
124 later are passed through a pooling layer that classifies the region as text or  
125 not text. In the model presented by [40], rotated region proposals are pre-  
126 sented, mostly to handle arbitrary oriented text. Analogously, R2CNN [41]  
127 the Region of Interest(ROI) pooling stage uses different fixed sizes which are  
128 concatenated for regression and classification. The work conducted by [42]  
129 mainly focuses on adaptive weighted pooling in different scales to further  
130 predict and regress region proposals.

## 131 2.2. Scene Text Recognition

132 Initial approaches explored by Jaderberg *et al.* [43] tackle text recognition  
133 as a classification problem. After training a CNN on synthetic generated  
134 samples, the obtained features are used to predict a vector that classifies  
135 the input word over approximately 90,000 classes. After the introduction  
136 of the Connectionist Temporal Classification (CTC) by Graves *et al.* [44]  
137 in handwriting recognition, the same methodology has been widely used in  
138 scene text as well. The work proposed by [45] employs the CTC layer after  
139 passing the input image through a CNN that acts as the encoder and a RNN

140 that act as the decoder. The introduction of an attention mechanism was  
141 initially proposed by [46] in the task of machine translation. This mechanism  
142 was briefly adopted in several vision tasks, including text recognition. The  
143 work proposed by [47], namely Focus Attention Network, employs attention  
144 to supervise relevant locations for word recognition. Bai *et al.* [48] introduce  
145 an edit probability to handle the misalignment between the ground truth  
146 string and the attention output string. Jaderberg *et al.* [49] proposed the  
147 Spatial Transformer Network, which is used by [50] to align detected text  
148 horizontally to further employ an attention based recognizer.

### 149 *2.3. End-to-End Text Recognition*

150 A commencing approach proposed by Jaderberg *et al.* [23] employs a  
151 sliding window to extract proposals, which are filtered and a CNN is used to  
152 regress the bounding boxes. Later the filtered regions that surpass a threshold  
153 are classified. In another work, Gupta *et al.* [51] defined a Fully Convolu-  
154 tional Regression Network for text detection and bounding box regression  
155 and the same classification network proposed by [23] for text recognition,  
156 being one of the first models that were fully trainable based on deep learning  
157 methodologies solely. In [52] a YOLO[53] based CNN is adopted to detect  
158 text instances, which later are passed through a Connectionist Temporal  
159 Classification module for recognition. These two stages are trained sepa-  
160 rately and later connected together to form an end-to-end architecture.  
161 The research presented by [54] introduces a CNN that is used as an encoder  
162 and a Long Short-Term Memory (LSTM) along with an attention mecha-  
163 nism module as decoder, both employed for detection and recognition. He *et*  
164 *al.* [55] use a CNN to extract proposals, which are fed into an LSTM to refine

165 the bounding boxes that are later employed as input to yet another LSTM  
166 to perform recognition that fixes misalignment between attention maps and  
167 ground truth character labels. In more recent work, [56] uses EAST [35] to  
168 obtain text regions and employs a CTC recognition module [44] to obtain  
169 an end-to-end reading system. Lyu *et al.* [57] use a variation of Mask R-  
170 CNN[39] to detect text in arbitrary shapes and segment an image in different  
171 instances to recognize similar text regions.

#### 172 2.4. Scene Text Retrieval

173 Closely related to our work, the scene text retrieval problem slightly dif-  
174 fers from classical scene text recognition methodologies. In a retrieval sce-  
175 nario the user defines a textual query which he wants to retrieve, whereas  
176 most of recognition approaches are based on employing a predefined vocab-  
177 ulary of the words one might come along within scene images. For instance,  
178 both Mishra *et al.*[22], who introduced the scene text retrieval task, and  
179 Jaderberg *et al.* [23], use a fixed vocabulary to create an inverted index  
180 which contains the presence of a word in the image. These approaches limit  
181 the freedom of queries to a set of predefined vocabulary words.

182 To address such a problem, text string descriptors based on n-gram frequen-  
183 cies, like the PHOC descriptor (Figure 2), have been successfully used for  
184 word spotting applications [58, 27, 59]. By using a vectorial codification of  
185 text strings, users can query any string at inference time without being lim-  
186 ited to a specific set of predefined vocabulary words. In this work, we make  
187 use of the PHOC descriptor along with an object detection framework based  
188 on YOLO [25, 53] that encodes found text instances. We suggest that this  
189 approach brings many benefits, mostly due to the high recall and single shot

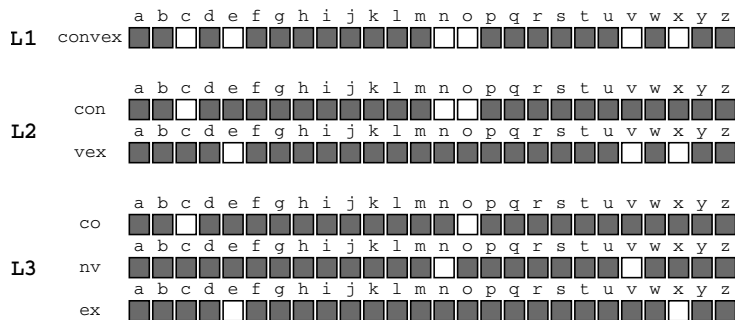


Figure 2: Pyramidal histogram of characters (PHOC) [27] of the word “convex” at levels 1, 2, and 3. The final PHOC representation is the concatenation of the partial one-hot encodings.

190 calculation required to locate and recognize text contained within an image,  
 191 accompanied by unmatched processing speeds.

## 192 2.5. Other applications

193 Fine-grained Classification is the task of classifying visually similar ob-  
 194 jects in which subtle differences are key to find discriminative features be-  
 195 tween classes. Finding these subtle features is a challenging task which keeps  
 196 this problem as an active topic in computer vision. Karaoglu *et al.* [20]  
 197 tackles this task by extracting visual features by employing a GoogleNet [60]  
 198 and a feature of Bag of Words to represent the text instances found in an  
 199 image and further classify them. More recently, [61] uses a similar approach  
 200 and extracts the visual features using a GoogleNet [60] and a combination  
 201 of two models: [32] to detect text and [45] to recognize text. The recognized  
 202 text instances are represented by GloVe [62], which are later used with an  
 203 attention mechanism on the visual features to classify the image.

204 Additional work has explored other fields of scene understanding by employ-  
 205 ing textual cues. The work proposed by [6] and [7] focuses on the Visual  
 206 Question Answering (VQA) [63] task. The VQA problem consists in provid-

207 ing an answer to a given image and question presented in natural language.  
208 Providing the correct answer is possible only if the system is capable to  
209 leverage textual information contained in the image.

### 210 **3. Proposed Architecture**

211 The proposed architecture is based on a custom-built YOLOv2 object  
212 detection model introduced by [25, 26]. This work adapts the object detector  
213 to output a compact representation of the text instances and recast them as  
214 a PHOC [27], thus enforcing the model to learn to construct such a vectorial  
215 codification. The suggested model is kept as a Fully Convolutional Neural  
216 Network, and a straightforward diagram is illustrated in Figure 3.

217 The convolutional neural network is composed of 22 convolutional layers  
218 with a leaky ReLU activation function after each convolution operation. The  
219 details of the proposed architecture can be seen in Table 1.

220 Batch normalization is used after every convolutional layer to help the  
221 model reach convergence. In total the model employs 5 max pooling layers,  
222 which reduces the input width and height by a factor of  $2^5$ . The filter size  
223 used in convolutions is  $3 \times 3$  and the channel number is doubled after each  
224 pooling step as in previous works that adopt a VGG [64] model backbone  
225 such as the work presented by [45]. In order to apply dimensionality reduc-  
226 tion and decrease the computation cost, the strategy proposed by the usage  
227 of an Inception module [60] is taken, and filters of size  $1 \times 1$  are interleaved  
228 between the  $3 \times 3$  convolutional filters to obtain richer feature maps. As  
229 it is defined in YoloV2 [26] and inspired in the Residual blocks introduced  
230 by [65], the convolutional backbone uses a pass-through layer from an earlier

Table 1: Detailed description of the proposed CNN architecture considering an input image size of 608 x 608.

Layer	Type	Filters	Size/Pad/Stride	Output
0	Input	-	-	608 x 608 x 3
1	Conv	32	3x3/p1/1	608 x 608 x 32
2	Max Pool		2x2/p0/2	304 x 304 x 32
3	Conv	64	3x3/p1/1	304 x 304 x 64
4	Max Pool		2x2/p0/2	152 x 152 x 64
5	Conv	128	3x3/p1/1	152 x 152 x 128
6	Conv	64	1x1/p0/1	152 x 152 x 64
7	Conv	128	3x3/p1/1	152 x 152 x 128
8	Max Pool		2x2/p0/2	76 x 76 x 128
9	Conv	256	3x3/p1/1	76 x 76 x 256
10	Conv	128	1x1/p0/1	76 x 76 x 128
11	Conv	256	3x3/p1/1	76 x 76 x 256
12	Max Pool		2x2/p0/2	38 x 38 x 256
13	Conv	512	3x3/p1/1	38 x 38 x 512
14	Conv	256	1x1/p0/1	38 x 38 x 256
15	Conv	512	3x3/p1/1	38 x 38 x 512
16	Conv	256	1x1/p0/1	38 x 38 x 256
17	Conv	512	3x3/p1/1	38 x 38 x 512
18	Max Pool		2x2/p0/2	19 x 19 x 512
19	Conv	1024	3x3/p1/1	19 x 19 x 1024
20	Conv	512	1x1/p0/1	19 x 19 x 512
21	Conv	1024	3x3/p1/1	19 x 19 x 1024
22	Conv	512	1x1/p0/1	19 x 19 x 512
23	Conv	1024	3x3/p1/1	19 x 19 x 1024
24	Conv	1024	3x3/p1/1	19 x 19 x 1024
26	Conv	1024	3x3/p1/1	19 x 19 x 1024
26	Concat[16]			38 x 38 x 512
27	Conv	64	1x1/p0/1	38 x 38 x 64
28	Concat[24,27]			19 x 19 x 1280
29	Conv	1024	3x3/p1/1	19 x 19 x 1024
30	Conv	7917	1x1/p0/1	19 x 19 x 7917

convolutional layer, which is concatenated and followed by a final  $1 \times 1$  convolutional filter with a linear activation with the number of filters matching the desired output tensor size to encode the PHOC descriptor.

Following the approach from the YOLOv2 model, we could define the word spotting task as a classification problem, where each detected word is a class. This one hot classification vector in the output tensor would represent the word class probability distribution among a defined list of words (fixed size dictionary) per each bounding box prediction. As simple as it sounds, such an approach limits the number of words that the model is able to recognize. In principle, if such a model requires to recognize 20 words, it would theoretically perform as well as classifying the 20 object classes from the PASCAL dataset presented in [26]. However, the problem raises in complex-



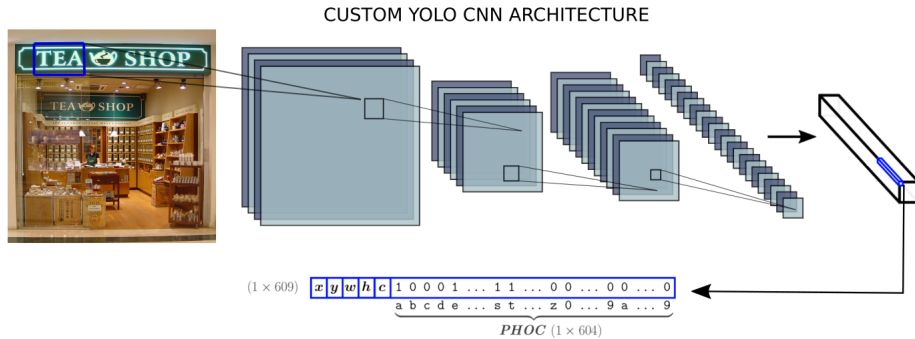


Figure 3: Our Convolutional Neural Network predicts at the same time bounding box coordinates  $x, y, w, h$ , an objectness score  $c$ , and a pyramidal histogram of characters (PHOC) of the word in each bounding box.

243 ity as the number of classes grow. If we consider training such a model (e.g.  
 244 the list of 90,000 most frequent words from the English vocabulary [23]),  
 245 the final convolutional layer would require 90,000 filters. This factor would  
 246 require an immense amount of data to successfully train such a model. Even  
 247 though a model with such characteristics could be designed, the limitation  
 248 of only recognizing words that belong to a predefined dictionary would still  
 249 be present. Recognizing out of vocabulary words would require a special  
 250 treatment or simply it would be a non-viable task. Furthermore, given the  
 251 number of parameters required, the model size would be too big and the real  
 252 time processing speed would most likely be lost.

253 A way of addressing the aforementioned problems, specifically a model that is  
 254 able to generalize and recognize previously unseen words, is desired. This is  
 255 the main driving rationale behind casting the network as a PHOC predictor,  
 256 which also permits to decrease the model's last filter size, thus allowing it  
 257 to perform at real-time. The PHOC [27] descriptor is a multi-level vectorial  
 258 representation of text strings that focuses on encoding if a specific character

259 is present in a defined spatial region of a string (see Figure 2). Intuitively,  
260 a CNN based model that effectively learns to predict the PHOC represen-  
261 tation of a detected word will inherently learn to identify the existence of  
262 a specific character in a visual region of the proposed bounding box. The  
263 model therefore will learn to construct the PHOC by automatically learn-  
264 ing character attributes independently. Learning how to construct such a  
265 representation given the morphology of a string allows the proposed model  
266 to transfer knowledge acquired at training time and employ it at inference  
267 time to build PHOCs of unseen words. This effect is possible due to the fact  
268 that the presence of a character at a particular locality of the word trans-  
269 lates to the same information in the PHOC representation, independently of  
270 the positioning or existence of other characters in the word. Moreover, the  
271 PHOC descriptor acts as a universal encoding scheme that offers unlimited  
272 expressiveness as it can represent any word constrained only by a language  
273 specific alphabet.

274 The PHOC version we propose in this work, contains a fixed length of  
275 604 dimensions represented as a binary vector.

276 In order to adapt the YOLOv2 object detection network for single shot  
277 detection and PHOC prediction, it is necessary to define the nature of the  
278 proposed descriptor. In the first place, the PHOC descriptor does not resem-  
279 ble a one hot vector as in a classification scheme. To treat the PHOC as a  
280 multi-hot binary vector, the last layer does not employ a softmax function.  
281 Secondly, the prediction of a PHOC vector is comprised of a set of numbers  
282 that satisfy the condition given by:

$$S = \{x|x \in \mathbb{R}, 0 \leq x \leq 1\} \quad (1)$$

283 Where  $S$  represents the set of possible PHOC values. In order to have  
 284 such a representation, a sigmoid activation function after the last convolu-  
 285 tional layer is used to predict the PHOC vectors rather than the original  
 286 softmax function.

287 Thirdly, we modify the original YOLOv2 Loss Function to facilitate the con-  
 288 vergence and learning process of the model. As it is presented in the original  
 289 YOLOv2 paper, the proposed algorithm is trained with the following multi-  
 290 part loss function:

$$L(b, C, c, \hat{b}, \hat{C}, \hat{c}) = \lambda_{box}L_{box}(b, \hat{b}) + L_{obj}(C, \hat{C}, \lambda_{obj}, \lambda_{noobj}) + \lambda_{cls}L_{cls}(c, \hat{c}) \quad (2)$$

291 where  $b$  is a vector with coordinates' offsets to an anchor bounding box,  $C$  is  
 292 the probability of that bounding box containing an object,  $c$  is the one hot  
 293 classification vector, and the three terms  $L_{box}$ ,  $L_{obj}$ , and  $L_{cls}$  are respectively  
 294 independent losses for bounding box regression, objectness estimation, and  
 295 classification. All the aforementioned losses are essentially the sum-squared  
 296 errors of ground truth  $(b, C, c)$  and predicted  $(\hat{b}, \hat{C}, \hat{c})$  values. At the moment  
 297 of predicting a PHOC,  $c$  (the ground truth) is a binary vector and  $\hat{c}$  (pre-  
 298 diction) meets the condition stated in 1, reason to opt for cross-entropy loss  
 299 function in  $L_{cls}$  as in a multi-label classification task:

$$L_{cls}(c, \hat{c}) = c \log \hat{c} + (1 - c) \log(1 - \hat{c}) \quad (3)$$

300 It is important to note that the combination of the sum-squared errors

301  $L_{box}$  and  $L_{obj}$  with the cross-entropy loss  $L_{cls}$  is controlled by the scaling  
 302 parameters  $\lambda_{box}$ ,  $\lambda_{obj}$ ,  $\lambda_{noobj}$ , and  $\lambda_{cls}$ .

303 Apart from the modifications made so far on top of the original YOLOv2  
 304 architecture we also changed the number, the scales, and the aspect ratios  
 305 of the pre-defined anchor boxes used by the network to predict bounding  
 306 boxes. Similar to [25], we have found the ideal set of anchor boxes  $B$  for our  
 307 training dataset by requiring that for each bounding box annotation there  
 308 exists at least one anchor box in  $B$  with an intersection over union of at least  
 309 0.6. Figure 4 illustrates the 13 bounding boxes found to be better suited for  
 310 our training data and their difference with the ones used in object detection  
 311 models.

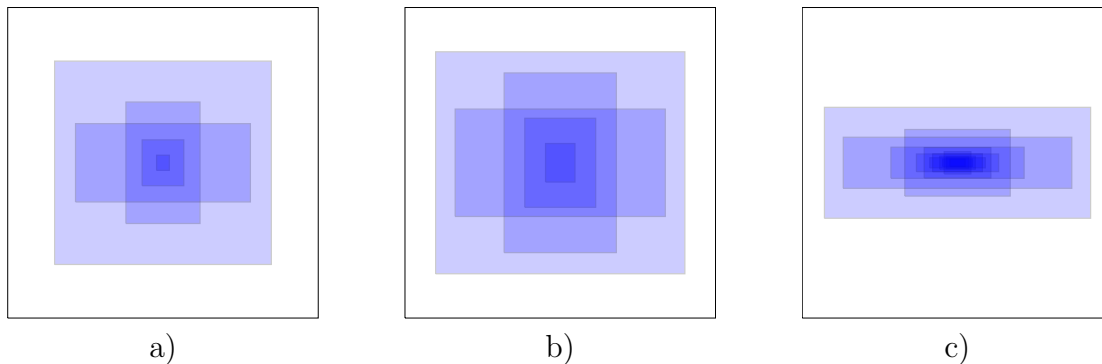


Figure 4: Anchor boxes used in the original YOLOv2 model for object detection in COCO (a) and PASCAL (b) datasets. (c) Our set of anchor boxes for text detection.

312 At test time, our model provides a total of  $W/32 \times H/32 \times 13$  bound-  
 313 ing box proposals, with  $W$  and  $H$  being the image input size, each one of  
 314 them with an objectness score ( $\hat{C}$ ) and a PHOC prediction ( $\hat{c}$ ). The original  
 315 YOLOv2 model filters the bounding box candidates with a detection thresh-  
 316 old  $\tau$  considering that a bounding box is a valid detection if  $\hat{C}max(\hat{c}) \geq \tau$ . If

317 the threshold condition is met, a non-maximal suppression (NMS) strategy  
318 is applied in order to get rid of overlapping detections of the same object. In  
319 our case the threshold is applied only on the objectness score ( $\hat{C}$ ) but with  
320 a much smaller value ( $\tau = 0.0025$ ) than in the original model ( $\tau \approx 0.2$ ), and  
321 we do not apply NMS. The reason is that any evidence of the presence of a  
322 word, even if it is small, it may be beneficial in terms of retrieval if its PHOC  
323 representation has a small distance to the PHOC of the queried word. With  
324 this threshold we generate an average of 50 descriptors for every image in  
325 the dataset and all of them form our retrieval database.

326 In this way, the scene text retrieval of a given query word is performed  
327 with a simple nearest neighbor search of the query PHOC representation over  
328 the outputs of the CNN in the entire image database. While the distance  
329 between PHOCs is usually computed using the cosine similarity, we did not  
330 find any noticeable downside on using an Euclidean distance for the nearest  
331 neighbor search.

### 332 *3.1. Training details*

333 We have trained our model in a modified version of the synthetic dataset  
334 of Gupta *et al.*[51]. First the dataset generator has been evenly modified  
335 to use a custom dictionary with the 90K most frequent English words, as  
336 proposed by Jaderberg *et al.*[23], instead of the Newsgroup20 dataset [66]  
337 dictionary originally used by Gupta *et al.*. The rationale was that in the  
338 original dataset there was no control over the word occurrences, and the  
339 distribution of word instances had a large bias towards stop-words found in  
340 newsgroups' emails. Moreover, the text corpus of the Newsgroup20 dataset  
341 contains words with special characters and non ASCII strings that we do



Figure 5: Synthetic training data generated with a modified version of the method of Gupta *et al.* [51]. We make use of a custom dictionary with the 90K most frequent English words, and restrict the range of random rotation to 15 degrees.

342 not contemplate in our PHOC representations. Finally, since the PHOC  
 343 representation of a word with a strong rotation does not make sense under  
 344 the pyramidal scheme employed, the dataset generator was modified to allow  
 345 rotated text up to 15 degrees. This way we generated a dataset of 1 million  
 346 images for training purposes. Figure 5 shows a set of samples of our training  
 347 data.

348 The model was trained for 30 epochs of the dataset using SGD with  
 349 a batch size of 64, an initial learning rate of 0.001, a momentum of 0.9,  
 350 and a decay of 0.0005. We initialize the weights of our model with the  
 351 YOLOv2 backbone pre-trained on Imagenet. During the firsts 10 epochs  
 352 we train the model only for word detection, without backpropagating the  
 353 loss of the PHOC prediction and using a fixed input size of  $448 \times 448$ . On  
 354 the following 10 epochs we start learning the PHOC prediction output with  
 355 the  $\lambda_{cls}$  parameter set to 1.0. After that, we continue learning for 10 more  
 356 epochs with a learning rate of 0.0001 and setting the parameters  $\lambda_{box}$  and  $\lambda_{cls}$   
 357 to 5.0 and 0.015 respectively. At this point we also adopted a multi-resolution  
 358 training, by randomly resizing the input images among 14 possible sizes in the  
 359 range from  $352 \times 352$  to  $800 \times 800$ , and we added new samples in our training

360 data. In particular, the added samples were the 1,233 training images of the  
361 ICDAR2013 [24] and ICDAR2015 [13] datasets. During the whole training  
362 process we used the same basic data augmentation as proposed by [25].

## 363 4. Experiments and results

364 In this section we present the experiments and results obtained on differ-  
365 ent standard benchmarks for text based image retrieval. First, we describe  
366 the datasets used throughout our experiments and after that, we present  
367 our results and compare them with the published state-of-the-art. As an  
368 extension to our previous work [29], an assessment when dealing with out-of-  
369 vocabulary words is conducted by analyzing the model in two multi-lingual  
370 datasets. Additionally, we conduct robustness experiments when confronted  
371 with imperfect image conditions, which further shows our models' poten-  
372 tial. Finally, we present a real-time text spotting application in videos, only  
373 possible by the characteristic speed capability of our method.

### 374 4.1. Datasets

#### 375 4.1.1. IIIT Scene Text Retrieval (STR)

376 The STR dataset [22] is a scene text image retrieval dataset composed  
377 of 10,000 images collected from the Google image search engine and Flickr.  
378 The dataset has 50 predefined query words and for each of them a list of  
379 10 – 50 relevant images (that contain the query word) is provided. It is  
380 a challenging dataset where relevant text appears in many different fonts  
381 and styles, and from different view points, among many distractors (images  
382 without any text).

383 *4.1.2. IIT Sports-10k dataset*

384 The Sports-10k dataset [22] is another scene text retrieval dataset com-  
385 posed of 10,000 images extracted from sports video clips. It has 10 pre-  
386 defined query words with their corresponding relevant images' lists. Scene  
387 text retrieval in this dataset is specially challenging because images are low  
388 resolution and often noisy or blurred, with small text generally located on  
389 advertisements signboards.

390 *4.1.3. Street View Text (SVT) dataset*

391 The SVT dataset [67] is comprised of images harvested from Google Street  
392 View where advertisement signboards is present. It contains more than 900  
393 words annotated in 350 different images. In our experiments we use the  
394 official partition that splits the images in a train set of 100 images and a  
395 test set of 249 images. This dataset also provides a lexicon of 50 words per  
396 image for recognition purposes, but we do not make use of it. For the image  
397 retrieval task we consider as queries the 427 unique words annotated on the  
398 test set.

399 *4.1.4. Multi-lingual scene text (MLT) datasets*

400 These two datasets MLT2017 [68] and MLT2019 [69] are scene text de-  
401 tection and recognition datasets that contain 7,200 and 10,000 images re-  
402 spectively in 10 different languages (Chinese, Japanese, Korean, English,  
403 French, Arabic, Italian, German, Bangla and Hindi) in equal proportions,  
404 representing 7 different scripts. These datasets mostly comprises focused  
405 text in natural images, and even though the main task is text detection and  
406 recognition, we adapted it to conduct text retrieval experiments. We employ  
407 this dataset to assess the generalization power of the PHOC representation



408 of unseen words at training time.

#### 409 4.1.5. *Text in videos (TiV) dataset*

410 The TiV dataset [70] contains 25 videos (13450 frames in total) and a test  
411 set of 24 videos (14374 frames in total) recorded from 4 different cameras.  
412 We use this dataset to assess the performance at real-time of our model at  
413 the moment of retrieving a specific text query. The challenge in this dataset  
414 remains in the fact that usually video frames contain a lower quality when  
415 compared to static images. The problems of text spotting usually relate to  
416 rotation, blur and occlusion of text found on each frame due to movement  
417 and focusing issues while including loss of information at the moment of video  
418 compression.

#### 419 4.2. *Scene text retrieval*

420 In the scene text retrieval task, the goal is to retrieve all images that con-  
421 tain instances of the query words in a dataset partition. Given a query, the  
422 database elements are sorted with respect to the probability of containing  
423 the queried word. We use the mean average precision as the accuracy mea-  
424 sure, which is the standard measure of performance for retrieval tasks and  
425 is essentially equivalent to the area below the precision-recall curve. Notice  
426 that, since the system always returns a ranked list with all the images in the  
427 dataset, the recall is always 100%. An alternative performance measure con-  
428 sists in considering only the top- $n$  ranked images and calculating the precision  
429 at this specific cut-off point ( $P@n$ ).

430 Table 2 compares the proposed method to previous state of the art for  
431 text based image retrieval on the IIIT-STR, Sports-10K, and SVT datasets.  
432 We show the mean average precision (mAP) and processing speed for the

433 same trained model using two different input sizes ( $576 \times 576$  and  $608 \times 608$ ),  
 434 and a multi-resolution version that combines the outputs of the model at  
 435 three resolutions (544, 576 and 608). Processing time has been calculated  
 436 using a Titan X (Pascal) GPU with a batch size of 1. We appreciate that  
 437 our method clearly outperforms previously published methods in two of the  
 438 benchmarks while it shows a competitive performance on the SVT dataset. It  
 439 is important to witness that our method achieves the highest measurements  
 440 in frames per second (fps), leading to the the best overall trade-off between  
 441 performance and processing speed in all datasets. Table 3 further compares  
 442 the proposed method to previous state of the art by showcasing the precision  
 443 at 10 (P@10) and 20 (P@20) on the Sports-10K dataset.

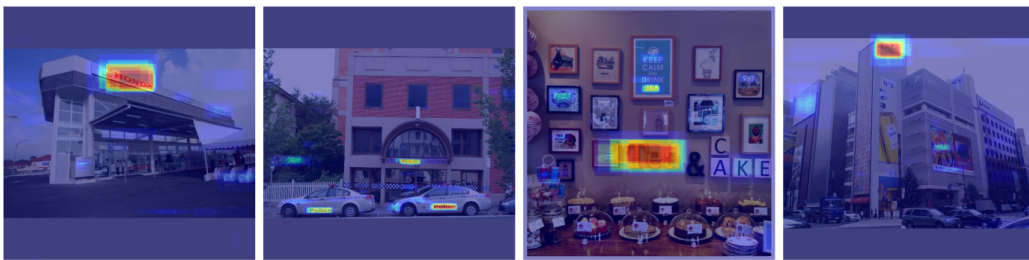


Figure 6: Bounding box heat-maps for queried words "honda", "police", "tea" and "sony" respectively.

444 In Figure 6, we depict the heat-maps of our model by calculating the  
 445 closests matching PHOC and its bounding box in relation to a given query.  
 446 As it can be seen on the showcased figure, several predicted PHOCs closely  
 447 match the queried word. Considering the implementation details defined  
 448 in the previous section, we avoid using a NMS post processing strategy to  
 449 preserve high matching PHOC proposals that could be discarded otherwise.

450 For a further analysis of the errors made by our model we have manually

Table 2: Comparison to previous state of the art for text based image retrieval: mean average precision (mAP) for IIIT-STR, and Sports-10K, and SVT datasets. (\*) Results reported by Mishra et al. in [22], not by the original authors. (†) Results computed with publicly available code from the original authors.

Method	STR (mAP)	Sports (mAP)	SVT (mAP)	fps
SWT [71]+ Mishra et al. [72]	-	-	19.25	
Wang <i>et al.</i> [67]	-	-	21.25*	
TextSpotter [73]	-	-	23.32*	1.0
Mishra <i>et al.</i> [22]	42.7	-	56.24	0.1
Ghosh <i>et al.</i> [74]	-	-	60.91	
Mishra [75]	44.5	-	62.15	0.1
Almazán <i>et al.</i> [27]	-	-	79.65	
TextProposals [76] + DictNet [43]	64.9†	67.5†	85.90†	0.4
Jaderberg <i>et al.</i> [23]	66.5	66.1	<b>86.30</b>	0.3
Bušta <i>et al.</i> [77] ICCV 2017	62.94	59.62	69.37	44.21
He <i>et al.</i> [55] CVPR 2018	50.16	50.74	72.82	1.25
He <i>et al.</i> [55] (With dictionary)	66.95	74.27	80.54	2.35
He <i>et al.</i> [55] (PHOC)	46.34	52.04	57.61	2.35
Proposed ( $576 \times 576$ )	<b>68.13</b>	<b>72.99</b>	82.02	<b>53.0</b>
Proposed ( $608 \times 608$ )	<b>69.83</b>	<b>73.75</b>	83.74	43.5
Proposed (multi-res.)	<b>71.37</b>	<b>74.67</b>	85.18	16.1

451 inspected the output of our model as well as the ground truth for the five  
452 queries with a lower mAP on the IIIT-STR dataset: "ibm", "indian", "insti-  
453 tute", "technology" and "sale". In most of these queries the low accuracy of  
454 our model can be explained in terms of having only very small and blurred  
455 instances in the database. In the case of "ibm", the characteristic font type  
456 in all instances of this word tends to be ignored by our model, and the  
457 same happens for some computer generated images (non scene images) that  
458 contain the word "sale". Figure 7 shows some examples of those instances.

Table 3: Comparison to previous state of the art for text based image retrieval: precision at n (P@n) for Sports-10K dataset.

Method	Sport-10K (P@10)	Sport-10K (P@20)
Mishra <i>et al.</i> [22]	44.82	43.42
Mishra [75]	47.20	46.25
Jaderberg <i>et al.</i> [23]	91.00	<b>92.50</b>
Proposed (576 × 576)	91.00	90.50
Proposed (multi-res.)	<b>92.00</b>	90.00



Figure 7: Error analysis: last ranked images for queries "sale", "ibm", "indian", "institute", "technology" and "police". Most of the errors made by our model come from text instances with a particular style, font type, size, etc. that is not well represented in our training data.

459 The analysis indicates that while our model is able to generalize well for  
 460 text strings not seen at training time it does not perform properly with text  
 461 styles, fonts, sizes not seen before. Our intuition is that this problem can be  
 462 alleviated with a richer training dataset.

#### 463 4.3. Multi-Lingual Scene Text Retrieval

464 As an extension to our previous work [29], we focus on analyzing the  
 465 generalization capability of the proposed model. It becomes essential to note  
 466 that designing an algorithm that learns to construct a compact representa-  
 467 tion of a string, such as the PHOC, paves the road to further development of  
 468 models that are not constrained to a fixed dictionary or training data sam-  
 469 ples. In order to assess the expressiveness of our architecture, we make use

470 of two Multi-lingual datasets 2017 [68] and 2019 [69] in which we can easily  
 471 find out-of-vocabulary words (text not seen at training time) with different  
 472 distributions and characteristics. These datasets are used by the research  
 473 community to perform text detection and recognition tasks, but not text  
 474 based image retrieval. Therefore, we have selected a set of 100 queries for in-  
 475 vocabulary experiments and another set of 100 queries for out-of-vocabulary  
 476 experiments for each dataset taken from the training split. Out-of-vocabulary  
 477 queries are selected by choosing the latin words with most occurrences af-  
 478 ter removing stop-words and words that contain non-alphanumeric charac-  
 479 ters. For in-vocabulary queries, we also remove stop-words and words with  
 480 non-alphanumeric characters before searching for latin words with similar  
 481 frequencies to the out-of-vocabulary queries.

Table 4: Comparison to previous state of the art method for text based image retrieval methods when queries are words already seen during the training process (IV) or not (OOV): mean average precision (mAP)

Method	MLT 2017		MLT 2019	
	IV	OOV	IV	OOV
He et al. [55]	24.79	19.47	27.6	24.99
Proposed	<b>46.52</b>	<b>46.87</b>	<b>46.41</b>	<b>46.03</b>

Table 5: Comparison to previous state of the art method for text based image retrieval methods when queries are words already seen during the training process (IV) or not (OOV): precision at n (P@n)

Method	MLT 2017						MLT 2019					
	IV			OOV			IV			OOV		
	P@5	P@10	P@20	P@5	P@10	P@20	P@5	P@10	P@20	P@5	P@10	P@20
He et al.	0.51	0.37	0.22	0.46	0.33	0.20	0.62	0.44	0.27	0.60	0.40	0.23
Proposed	<b>0.77</b>	<b>0.57</b>	<b>0.34</b>	<b>0.78</b>	<b>0.59</b>	<b>0.34</b>	<b>0.80</b>	<b>0.64</b>	<b>0.41</b>	<b>0.80</b>	<b>0.64</b>	<b>0.40</b>

482 Tables 4 and 5 show the ability for our model to perform retrieval with

483 the same accuracy for in-vocabulary queries and out-of-vocabulary queries in  
 484 both datasets. As we stated previously, this is because our model is learning  
 485 how to build a PHOC from text rather than performing a classification along  
 486 a fixed dictionary. It is important to note that our model performs signifi-  
 487 cantly better than a state of the art reading system presented by [55] at the  
 488 text retrieval task. Additionally, the method from [55] was trained using the  
 489 dictionary from [66] which contains English words, thus performing poorly  
 490 when dealing with out of vocabulary words mostly belonging to different lan-  
 491 guages. Figure 8 shows the top-5 ranked images for the queries "vodafone"  
 492 in IIIT-STR dataset, "uscita" (italian) in MLT 2017 and "werden" (german)  
 493 in MLT 2019, all of them being unseen samples at training time. In all of  
 494 them our model reaches a 100% precision at 5.



Figure 8: From top to bottom, top-5 ranked images for the queries "vodafone", "uscita", "werden". Although our model has not seen these words at training time it is able to achieve a 100% P@5 for all of them.

495 *4.4. Robustness of the Model*

496 In the following subsection, experiments to determine the robustness of  
497 the model to imperfect conditions are performed. Experiments regarding  
498 rotation, blur, compression and occlusion are analyzed.

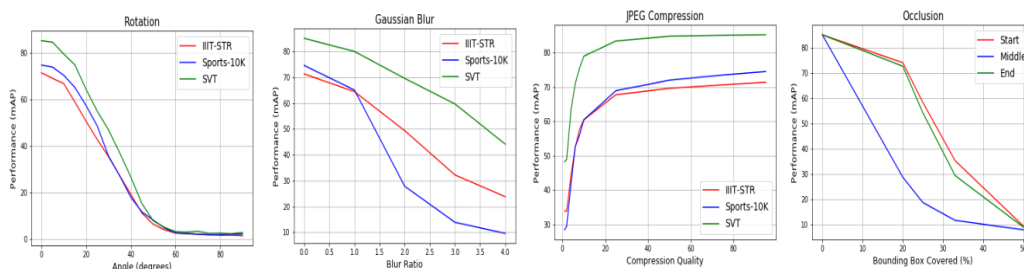


Figure 9: Robustness performance for imperfect conditions such as rotation, blur, compression and occlusion.

499 *4.4.1. Rotation*

500 A big difference between text found in documents and text in natural  
501 imagery is the arbitrary orientation text may have. Rotated and arbitrary  
502 shaped text instances are one of the main problems in the research com-  
503 munity. Challenges such as the one presented in [78] remains as an open  
504 problem and an active field, in which the task is far from being a trivial one.  
505 Experiments to assess the model performance and robustness towards rota-  
506 tion were conducted. Each image from the analyzed datasets was rotated by  
507 a specific angle starting at 0° to 90° in steps of 5°, clockwise and counter  
508 clockwise. The images were rotated by considering the center of the image  
509 as the reference point as it is show in Figure 10. Bi-linear interpolation was  
510 used in order to avoid losing information and padding was used in order to  
511 avoid cutting-off sections that contain text in an image.



Figure 10: From left to right, qualitative rotated sample image at  $0^\circ$ ,  $20^\circ$ ,  $40^\circ$ ,  $60^\circ$  and  $80^\circ$  taken from SVT dataset. Spatial positioning of characters is lost at high rotation angles, thus decreasing the model capability of constructing the PHOC representation.

512 As it can be seen in Figure 9, the greater the rotation angle is applied  
 513 in the image, the performance of the model decreases. The rotation effect is  
 514 amplified in the IIIT-STR dataset due to the fact that it already contains  
 515 text in different orientations when compared to the more stable and horizon-  
 516 tal text occurrences found in the remaining two datasets. It is worth noting  
 517 that the proposed model was trained employing a synthetic dataset that in-  
 518 cluded rotated words up to an angle of 15 degrees. This effect is perceived by  
 519 noticing a significant decrease in performance (increase in gradient) when-  
 520 ever an image is rotated more than 25 degrees. Another fact that decreases  
 521 the rotation performance in angles that approach to  $90^\circ$  is the shape of the  
 522 predefined anchor boxes (Figure 4 c.), which possess a shape that mostly  
 523 captures horizontal text. The orientation of text is key at the moment of  
 524 building the PHOC representation of a word. This representation is con-  
 525 structed by considering spatial information of each character contained in a  
 526 string, which is heavily affected by rotated words contained in an image.

#### 527 4.4.2. *Blur*

528 Blur in text is a common issue in incidental images [13] as well as in video  
 529 frames, specially on videos that contain rapid camera movement, fast scene  
 530 transitions and not professional cameras. Different kernel sizes of Gaussian



531 blur are employed to assess the proposed model performance, implemented by  
 532 using [79]. As it can be seen in the qualitative results depicted in Figure 11,  
 533 humans will not have a difficult time recognizing most of the text occurrences  
 534 in blurred images. Blur is a particularly big problem in the Sports-10K  
 535 dataset, which due to its nature, video frames depicted already contain blurry  
 536 text. Gaussian blur augments this issue, thus a sharp decrease in performance  
 537 is noted when compared to the remaining datasets, see Figure 9. Further  
 538 strategies of data augmentation with blurred images or de-blurring techniques  
 539 as presented by [80] can be used as an additional step before inference time.



Figure 11: Increasing Gaussian blur in a sample image taken from IIIT-STR dataset. Fine features that differentiate characters are lost, thus affecting the ability of the proposed model to recognize a word.

#### 540 4.4.3. Compression

541 Compression in images and video can severely degrade the image qual-  
 542 ity, thus affecting subtle details that impact the performance of deep nets.  
 543 In order to simulate real life compression issues, different lossy compression  
 544 qualities were employed to downgrade images in the proposed datasets by us-  
 545 ing the JPEG compression algorithm. The perception in quality degradation  
 546 is not linear, thus more emphasis was placed in extreme scenarios (low com-  
 547 pression qualities). The compression method used was taken from the public  
 548 implementation from [79] and different quality values were employed. As it  
 549 can be seen in the qualitative results depicted in Figure 12, changes in quality

550 above the value of 25 are barely noticeable by human perception alone. De-  
 551 spite this fact, as it can be seen in Figure 9, our model achieves a comparable  
 552 performance with previous state of the art methods depicted in Table 2 even  
 553 when the input image belongs to a low quality compression range. It is worth  
 554 pointing out that for quality values of 20 and above the performance gradient  
 555 tends to decrease making the performance grow slowly until achieving state  
 556 of the art reported values in images with a higher compression quality. Sim-  
 557 ilarly to the blur problems encountered in previous section, the Sports-10K  
 558 dataset is the most susceptible to low image qualities, due to the collecting  
 559 process of this dataset at the moment of extracting frames from video.

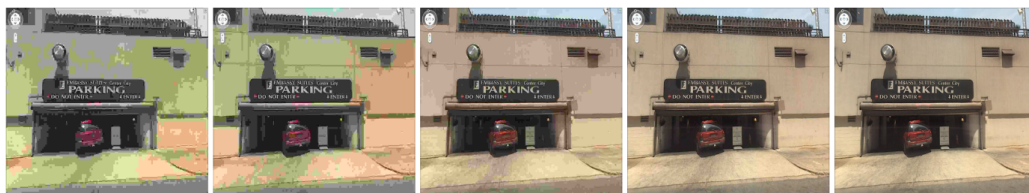


Figure 12: Increasing compression quality from left to right (1, 4, 8, 25, 75) in sample image from SVT dataset. At low qualities text at small scales is barely legible. Despite this effect, our model achieves state of the art level performance at qualities bigger than 20.

#### 560 4.4.4. Occlusion

561 An ongoing challenge in the scene text reading community is occlusion, as  
 562 it can entirely modify the morphology of spotted text. Humans are less prone  
 563 to occlusion problems, due to prior knowledge of the context of an image or  
 564 by the existing familiarity towards a specific language. In our experiments,  
 565 three scenarios were proposed according to the position of the occlusion,  
 566 namely at the beginning, middle and end of a word. These experiments were  
 567 conducted only in the SVT dataset because it was the only one that already

568 contained bounding box labels. The occlusion was generated by extreme  
569 blurring of a given percentage of the area that contains text in an image.  
570 The percentages of occlusion employed were half, one third, one fourth and  
571 one fifth of the total bounding box area, some qualitative samples can be  
572 seen in Figure 13. Not all text occurrences in a given image are occluded,  
573 because there are words that do not contain any ground truth annotations  
574 provided in the SVT dataset. As it can be seen in Figure 9, when the  
575 occlusion is located at the beginning and end of a word the model achieves  
576 a similar performance which slowly decreases as the occluded area grows.  
577 The model learns to build the PHOC of the occluded word, and successfully  
578 retrieves the closest matching representation. This outcome can be seen  
579 in Figure 14, in which the model successfully retrieves occluded images for  
580 the query "adidas". However, when the occlusion affects the center of a  
581 word, the model achieves a lower performance at the moment of retrieving  
582 a specific query. This outcome can be easily explained because the detected  
583 text is treated as two different word occurrences, thus generating different  
584 proposals that actually belong to the same word.



Figure 13: Occlusion samples. From left to right: occlusion located at the beginning of the image occupying 1/3 of the total bounding box area, occlusion at the beginning involving 1/5, occlusion at the middle filling 1/3, and occlusion at the end covering 1/3 and 1/5 of the total bounding boxes area respectively.



Figure 14: Images within the top 10 ranked images for the query "adidas". Our model successfully retrieves partially occluded and blurred words.

#### 585 4.5. Real-time Text Spotting in Videos

586 Given the high processing frame-rates that we achieve (c.f. Table 2),  
 587 we can use the proposed method for spotting text in video streams in real  
 588 time. Such application might be interesting in scenarios like assistance to  
 589 driving systems, in order to spot certain words in the open world, or to track  
 590 advertisement exposure in sports broadcasting. In such cases, the user casts  
 591 a textual query that has to be sought within videos. We shall take into  
 592 account that video recorded in natural scenes contain text instances that  
 593 are extremely susceptible to imperfect conditions. Low quality of recording  
 594 devices and rapid camera movement tends to produce blurred and rotated  
 595 content. Text found in video is also vulnerable to unintended occlusions that  
 596 affect several consecutive frames. In order to test the performance of the  
 597 proposed method in such scenario, we have used the Text in Videos challenge  
 598 dataset [24], in which the train partition consists of 25 videos, 13.450 frames  
 599 in total, with their corresponding ground-truth annotation. We decided to  
 600 use as queries the 20 words having more than three letters that have more  
 601 occurrences in the dataset. Having set a threshold on the distance between  
 602 the query PHOC representation and the closest word hypothesis in each  
 603 frame, we decide whether the queried word appears or not in that frame.  
 604 We evaluate the text spotting in videos task by using the F-score, so that we

Table 6: Top 15 most frequent words with their number of occurrences and the reached F-score.

Query	Occurrences	F-score
<i>flor</i>	539	94.05
<i>Marie</i>	426	83.89
<i>Renfe</i>	314	78.26
<i>createurs</i>	303	72.40
<i>Dixan</i>	278	87.54
<i>FONTANEDA</i>	261	84.44
<i>VOTRE</i>	257	91.01
<i>Digestive</i>	254	90.00
<i>USHIP</i>	245	75.35
<i>ACCASTILLEUR</i>	241	66.26
<i>Applus</i>	237	91.96
<i>Rectorat</i>	237	88.96
<i>CONSEIL</i>	230	83.18
<i>mundi</i>	230	85.24
<i>Accastillage</i>	199	61.41
<i>MISTOL</i>	186	57.51
<b>Average</b>	—	<b>76.70</b>

605 penalize both missing frames where the query word appears and false positive  
606 frames. Overall we achieved an F-score of 76.70, and we provide some results  
607 for the topmost 15 queries in Table 6. Video demos are available in our public  
608 repository<sup>1</sup>.

## 609 5. Conclusions

610 In this work, we presented a real-time performing word spotting method,  
611 based on a fully convolutional neural network that allows to detect and rec-  
612 ognize text in a single calculation which yields real-time processing capa-

<sup>1</sup><https://github.com/lluisgomez/single-shot-str>

613 bility. The introduced model significantly improves previous state of the  
614 art results on the scene text retrieval task on the IIIT-STR and Sports-10K  
615 dataset while obtaining comparable results to the state of the art in the SVT  
616 Dataset. Moreover, it can do so achieving speeds  $50\times$  to  $150\times$  speed com-  
617 pared to other state of the art methods, which opens up the possibility of  
618 employing this model for real time scenarios, such as video, and indexing  
619 large scale databases.

620 Importantly, it has been shown that the proposed method is able to con-  
621 struct a compact vectorial representation of out of dictionary queries at in-  
622 ference time, while keeping the performance at words previously seen at  
623 training. Achieving this result is possible by employing the PHOC as a word  
624 representation instead of tackling the task as a direct word classification.  
625 The method showcased is able to generalize unseen samples in a robust and  
626 efficient way, as the evidence strongly points out in experiments performed  
627 in a multilingual dataset. Additionally, the model proves to be robust at  
628 dealing with highly compressed images and text samples with occlusions at  
629 the beginning and at the end of a word. However, large rotation angles still  
630 present a problem which can be tackled by synthesizing training data with  
631 different characteristics and by using different priors when defining anchor  
632 boxes. Additional future work can be conducted to investigate the use of  
633 word embeddings that exploit the morphology of a word other than PHOC.

634 The code, pre-trained models, data and demo videos used in this work are  
635 publicly available at <https://github.com/lluisgomez/single-shot-str>.

636

637 **Funding**

638 This work was partially funded by the Spanish Research project TIN2017-  
639 89779-P, the grant given by the European Social Fund 2014-2020 (CCI:  
640 2014ES05SFOP007) and Predoctoral grant (AGAUR) number 2019-FIB01233,  
641 the H2020 MarieSkodowska-Curie actions of the European Union, grant agree-  
642 ment No 712949(TECNIOspring PLUS), and UAB PhD scholarship No B18P0070.

643 **6. Bibliography**

644 **References**

- 645 [1] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L.  
646 Zitnick, Microsoft COCO: Common objects in context, in: Proc. of the European  
647 Conference on Computer Vision, Springer, pp. 740–755.
- 648 [2] A. Veit, T. Matera, L. Neumann, J. Matas, S. Belongie, COCO-text: Dataset and  
649 benchmark for text detection and recognition in natural images, arXiv preprint  
650 arXiv:1601.07140 (2016).
- 651 [3] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, Nature 521 (2015) 436–444.
- 652 [4] S. S. Tsai, H. Chen, D. Chen, G. Schroth, R. Grzeszczuk, B. Girod, Mobile visual  
653 search on printed documents using text and low bit-rate features, in: 2011 18th IEEE  
654 International Conference on Image Processing, IEEE, pp. 2601–2604.
- 655 [5] G. Schroth, S. Hilsenbeck, R. Huitl, F. Schweiger, E. Steinbach, Exploiting text-  
656 related features for content-based image retrieval, in: 2011 IEEE International Sym-  
657 posium on Multimedia, IEEE, pp. 77–84.
- 658 [6] A. Singh, V. Natarajan, M. Shah, Y. Jiang, X. Chen, D. Batra, D. Parikh,  
659 M. Rohrbach, Towards vqa models that can read, in: The IEEE Conference on  
660 Computer Vision and Pattern Recognition (CVPR).
- 661 [7] D. Karatzas, Ll. Gómez, M. Rusiñol, A. Biten, A. Mafla, R. Tito, E. Valveny, C.  
662 Jawahar, M. Mathew, ICDAR 2019 Robust Reading Challenge on Scene Text Vi-  
663 sual Question Answering, <http://http://rrc.cvc.uab.es/?ch=11>, 2019. [Online,  
664 accessed 22-April-2019].
- 665 [8] Y. Dvorin, U. E. Havosha, Method and device for instant translation, 2009. US Patent  
666 App. 11/998,931.

- 667 [9] C. Parkinson, J. J. Jacobsen, D. B. Ferguson, S. A. Pombo, Instant translation sys-  
668 tem, 2016. US Patent 9,507,772.
- 669 [10] X. Liu, J. K. Samarabandu, A simple and fast text localization algorithm for indoor  
670 mobile robot navigation, in: *Image Processing: Algorithms and Systems IV*, volume  
671 5672, International Society for Optics and Photonics, pp. 139–151.
- 672 [11] X. Liu, J. Samarabandu, An edge-based text region extraction algorithm for in-  
673 door mobile robot navigation, in: *IEEE International Conference Mechatronics and*  
674 *Automation*, 2005, volume 2, IEEE, pp. 701–706.
- 675 [12] Y. Zhu, C. Yao, X. Bai, Scene text detection and recognition: Recent advances and  
676 future trends, *Frontiers of Computer Science* 10 (2016) 19–36.
- 677 [13] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura,  
678 J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, et al., ICDAR 2015 competition  
679 on robust reading, in: *Proc. of the IEEE International Conference on Document*  
680 *Analysis and Recognition*, pp. 1156–1160.
- 681 [14] Z. He, J. Liu, H. Ma, P. Li, A new automatic extraction method of container identity  
682 codes, *IEEE Transactions on intelligent transportation systems* 6 (2005) 72–78.
- 683 [15] M. A. Chowdhury, K. Deb, Extracting and segmenting container name from container  
684 images, *International Journal of Computer Applications* 74 (2013).
- 685 [16] I. Sutskever, G. E. Hinton, A. Krizhevsky, Imagenet classification with deep convo-  
686 lutional neural networks, *Advances in neural information processing systems* (2012)  
687 1097–1105.
- 688 [17] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpa-  
689 thy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge,  
690 *International journal of computer vision* 115 (2015) 211–252.
- 691 [18] S. Long, X. He, C. Ya, Scene text detection and recognition: The deep learning era,  
692 arXiv preprint arXiv:1811.04256 (2018).
- 693 [19] Y. Movshovitz-Attias, Q. Yu, M. C. Stumpe, V. Shet, S. Arnoud, L. Yatziv, Ontolog-  
694 ical supervision for fine grained classification of street view storefronts, in: *Proc. of*  
695 *the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1693–1702.
- 696 [20] S. Karaoglu, R. Tao, T. Gevers, A. W. Smeulders, Words matter: Scene text for image  
697 classification and retrieval, *IEEE Transactions on Multimedia* 19 (2017) 1063–1076.
- 698 [21] X. Bai, M. Yang, P. Lyu, Y. Xu, Integrating scene text and visual appearance for  
699 fine-grained image classification with convolutional neural networks, arXiv preprint  
700 arXiv:1704.04613 (2017).



- 701 [22] A. Mishra, K. Alahari, C. Jawahar, Image retrieval using textual cues, in: Proc. of  
702 the IEEE International Conference on Computer Vision, pp. 3040–3047.
- 703 [23] M. Jaderberg, K. Simonyan, A. Vedaldi, A. Zisserman, Reading text in the wild with  
704 convolutional neural networks, *International Journal of Computer Vision* 116 (2016)  
705 1–20.
- 706 [24] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i Bigorda, S. R. Mestre,  
707 J. Mas, D. F. Mota, J. A. Almazan, L. P. De Las Heras, ICDAR 2013 robust reading  
708 competition, in: Proc. of the IEEE International Conference on Document Analysis  
709 and Recognition, pp. 1484–1493.
- 710 [25] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-  
711 time object detection, in: Proc. of the IEEE Conference on Computer Vision and  
712 Pattern Recognition, pp. 779–788.
- 713 [26] J. Redmon, A. Farhadi, YOLO9000: better, faster, stronger, arXiv preprint  
714 arXiv:1612.08242 (2016).
- 715 [27] J. Almazán, A. Gordo, A. Fornés, E. Valveny, Word spotting and recognition with  
716 embedded attributes, *IEEE Transactions on Pattern Analysis and Machine Intelli-  
717 gence* 36 (2014) 2552–2566.
- 718 [28] S. Sudholt, G. A. Fink, Phocnet: A deep convolutional neural network for word  
719 spotting in handwritten documents, in: Proc. of the IEEE International Conference  
720 on Frontiers in Handwriting Recognition, pp. 277–282.
- 721 [29] L. Gómez, A. Mafla, M. Rusinol, D. Karatzas, Single shot scene text retrieval, in:  
722 Proceedings of the European Conference on Computer Vision (ECCV), pp. 700–715.
- 723 [30] D. He, X. Yang, C. Liang, Z. Zhou, A. G. Ororbi, D. Kifer, C. Lee Giles, Multi-scale  
724 fcn with cascaded instance aware segmentation for arbitrary oriented word spotting  
725 in the wild, in: Proceedings of the IEEE Conference on Computer Vision and Pattern  
726 Recognition, pp. 3519–3528.
- 727 [31] C. Yao, X. Bai, N. Sang, X. Zhou, S. Zhou, Z. Cao, Scene text detection via holistic,  
728 multi-channel prediction, arXiv preprint arXiv:1606.09002 (2016).
- 729 [32] M. Liao, B. Shi, X. Bai, X. Wang, W. Liu, Textboxes: A fast text detector with a sin-  
730 gle deep neural network, in: Proc. of the AAAI Conference on Artificial Intelligence,  
731 pp. 4161–4167.
- 732 [33] M. Liao, B. Shi, X. Bai, Textboxes++: A single-shot oriented scene text detector,  
733 arXiv preprint arXiv:1801.02765 (2018).
- 734 [34] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, SSD:  
735 Single shot multibox detector, in: Proc. of the European Conference on Computer  
736 Vision, Springer, pp. 21–37.

- 737 [35] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, J. Liang, East: an efficient  
738 and accurate scene text detector, in: 2017 IEEE Conference on Computer Vision and  
739 Pattern Recognition (CVPR), IEEE, pp. 2642–2651.
- 740 [36] K.-H. Kim, S. Hong, B. Roh, Y. Cheon, M. Park, Pvanet: deep but lightweight neural  
741 networks for real-time object detection, arXiv preprint arXiv:1608.08021 (2016).
- 742 [37] R. Girshick, Fast r-cnn, in: Proceedings of the IEEE international conference on  
743 computer vision, pp. 1440–1448.
- 744 [38] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards real-time object detection  
745 with region proposal networks, in: Proc. of the International Conference on Neural  
746 Information Processing Systems, pp. 91–99.
- 747 [39] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE  
748 international conference on computer vision, pp. 2961–2969.
- 749 [40] J. Ma, W. Shao, H. Ye, L. Wang, H. Wang, Y. Zheng, X. Xue, Arbitrary-oriented  
750 scene text detection via rotation proposals, IEEE Transactions on Multimedia 20  
751 (2018) 3111–3122.
- 752 [41] Y. Jiang, X. Zhu, X. Wang, S. Yang, W. Li, H. Wang, P. Fu, Z. Luo, R2cnn:  
753 Rotational region cnn for orientation robust scene text detection, arXiv preprint  
754 arXiv:1706.09579 (2017).
- 755 [42] S. Zhang, Y. Liu, L. Jin, C. Luo, Feature enhancement network: A refined scene text  
756 detector, in: Thirty-Second AAAI Conference on Artificial Intelligence.
- 757 [43] M. Jaderberg, K. Simonyan, A. Vedaldi, A. Zisserman, Synthetic data and artificial  
758 neural networks for natural scene text recognition, arXiv preprint arXiv:1406.2227  
759 (2014).
- 760 [44] A. Graves, S. Fernández, F. Gomez, J. Schmidhuber, Connectionist temporal clas-  
761 sification: labelling unsegmented sequence data with recurrent neural networks, in:  
762 Proceedings of the 23rd international conference on Machine learning, ACM, pp.  
763 369–376.
- 764 [45] B. Shi, X. Bai, C. Yao, An end-to-end trainable neural network for image-based  
765 sequence recognition and its application to scene text recognition, IEEE Transactions  
766 on Pattern Analysis and Machine Intelligence 39 (2017) 2298–2304.
- 767 [46] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to  
768 align and translate, arXiv preprint arXiv:1409.0473 (2014).
- 769 [47] Z. Cheng, F. Bai, Y. Xu, G. Zheng, S. Pu, S. Zhou, Focusing attention: Towards  
770 accurate text recognition in natural images, in: Proceedings of the IEEE International  
771 Conference on Computer Vision, pp. 5076–5084.

- 772 [48] F. Bai, Z. Cheng, Y. Niu, S. Pu, S. Zhou, Edit probability for scene text recog-  
773 nition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern  
774 Recognition, pp. 1508–1516.
- 775 [49] M. Jaderberg, K. Simonyan, A. Zisserman, et al., Spatial transformer networks, in:  
776 Advances in neural information processing systems, pp. 2017–2025.
- 777 [50] B. Shi, X. Wang, P. Lyu, C. Yao, X. Bai, Robust scene text recognition with au-  
778 tomatic rectification, in: Proc. of the IEEE Conference on Computer Vision and  
779 Pattern Recognition, pp. 4168–4176.
- 780 [51] A. Gupta, A. Vedaldi, A. Zisserman, Synthetic data for text localisation in nat-  
781 ural images, in: Proc. of the IEEE Conference on Computer Vision and Pattern  
782 Recognition, pp. 2315–2324.
- 783 [52] M. Busta, L. Neumann, J. Matas, Deep textspotter: An end-to-end trainable scene  
784 text localization and recognition framework, in: Proceedings of the IEEE Interna-  
785 tional Conference on Computer Vision, pp. 2204–2212.
- 786 [53] J. Redmon, A. Farhadi, Yolo9000: better, faster, stronger, arXiv preprint (2017).
- 787 [54] H. Li, P. Wang, C. Shen, Towards end-to-end text spotting with convolutional recur-  
788 rent neural networks, arXiv preprint arXiv:1707.03985 (2017).
- 789 [55] T. He, Z. Tian, W. Huang, C. Shen, Y. Qiao, C. Sun, An end-to-end textspotter  
790 with explicit alignment and attention, in: Proceedings of the IEEE Conference on  
791 Computer Vision and Pattern Recognition, pp. 5020–5029.
- 792 [56] X. Liu, D. Liang, S. Yan, D. Chen, Y. Qiao, J. Yan, Fots: Fast oriented text spotting  
793 with a unified network, arXiv preprint arXiv:1801.01671 (2018).
- 794 [57] P. Lyu, M. Liao, C. Yao, W. Wu, X. Bai, Mask textspotter: An end-to-end train-  
795 able neural network for spotting text with arbitrary shapes, in: Proceedings of the  
796 European Conference on Computer Vision (ECCV), pp. 67–83.
- 797 [58] D. Aldavert, M. Rusiñol, R. Toledo, J. Lladós, Integrating visual and textual cues  
798 for query-by-string word spotting, in: Proc. of the IEEE International Conference on  
799 Document Analysis and Recognition, pp. 511–515.
- 800 [59] S. K. Ghosh, E. Valveny, Query by string word spotting based on character bi-gram  
801 indexing, in: Proc. of the IEEE International Conference on Document Analysis and  
802 Recognition, pp. 881–885.
- 803 [60] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Van-  
804 houcke, A. Rabinovich, Going deeper with convolutions, in: Proc. of the IEEE  
805 Conference on Computer Vision and Pattern Recognition, pp. 1–9.

- 806 [61] X. Bai, M. Yang, P. Lyu, Y. Xu, J. Luo, Integrating scene text and visual appearance  
807 for fine-grained image classification, *IEEE Access* 6 (2018) 66322–66335.
- 808 [62] J. Pennington, R. Socher, C. Manning, Glove: Global vectors for word representation,  
809 in: *Proceedings of the 2014 conference on empirical methods in natural language  
810 processing (EMNLP)*, pp. 1532–1543.
- 811 [63] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, D. Parikh,  
812 Vqa: Visual question answering, in: *Proceedings of the IEEE international conference  
813 on computer vision*, pp. 2425–2433.
- 814 [64] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image  
815 recognition, *arXiv preprint arXiv:1409.1556* (2014).
- 816 [65] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in:  
817 *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.  
818 770–778.
- 819 [66] K. Lang, T. Mitchell, *Newsgroup 20 dataset* (1999).
- 820 [67] K. Wang, B. Babenko, S. Belongie, End-to-end scene text recognition, in: *Proc. of  
821 the IEEE International Conference on Computer Vision*, pp. 1457–1464.
- 822 [68] I. Bizid, J. Chazalon, H. Choi, Y. Feng, D. Karatzas, W. Khelif, Z. Luo, M. Luqman, N.  
823 Nayef, U. Pal, C. Rigaud, F. Yin J. Matas, N. Nayef, U. Pal, Y. Patel, *ICDAR2017  
824 Competition on Multi-lingual scene text detection and script identification*, <http://http://rrc.cvc.uab.es/?ch=8>, 2017. [Online, accessed 22-April-2019].
- 826 [69] M. Bušta, D. Karatzas, W. Khelif, J. Matas, N. Nayef, U. Pal, Y. Patel, *ICDAR 2019  
827 Robust Reading Challenge on Multi-lingual scene text detection and recognition*,  
828 <http://http://rrc.cvc.uab.es/?ch=11>, 2019. [Online, accessed 30-April-2019].
- 829 [70] M. Iwamura, L. Gomez, D. Karatzas, *Robust Reading Challenge on Text in videos  
830 2013-2015*, <http://http://rrc.cvc.uab.es/?ch=11>, 2019. [Online, accessed 22-  
831 April-2019].
- 832 [71] B. Epshtein, E. Ofek, Y. Wexler, Detecting text in natural scenes with stroke width  
833 transform, in: *Proc. of the IEEE Conference on Computer Vision and Pattern Recog-  
834 nition*, pp. 2963–2970.
- 835 [72] A. Mishra, K. Alahari, C. Jawahar, Top-down and bottom-up cues for scene text  
836 recognition, in: *Proc. of the IEEE Conference on Computer Vision and Pattern  
837 Recognition*, pp. 2687–2694.
- 838 [73] L. Neumann, J. Matas, Real-time scene text localization and recognition, in: *Proc. of  
839 the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3538–3545.

- 840 [74] S. K. Ghosh, L. Gomez, D. Karatzas, E. Valveny, Efficient indexing for query by  
841 string text retrieval, in: Proc. of the IEEE International Conference on Document  
842 Analysis and Recognition, pp. 1236–1240.
- 843 [75] A. Mishra, Understanding Text in Scene Images, Ph.D. thesis, International Institute  
844 of Information Technology Hyderabad, 2016.
- 845 [76] L. Gómez, D. Karatzas, Textproposals: a text-specific selective search algorithm for  
846 word spotting in the wild, Pattern Recognition 70 (2017) 60–74.
- 847 [77] M. Buvsta, L. Neumann, J. Matas, Deep textspotter: An end-to-end trainable scene  
848 text localization and recognition framework, in: Proc. of the IEEE International  
849 Conference on Computer Vision, pp. 2204–2212.
- 850 [78] Y. Sun, C Kheng, C. Chet, Y. Liu, C. Luo, Z. Ni, D. Karatzas, S. Zhang, J. Han, E.  
851 Ding, C. Seng, L. Jin, ICDAR2019 Robust Reading Challenge on Arbitrary-Shaped  
852 Text, <http://rrc.cvc.uab.es/?ch=14&com=introduction>, 2019. [Online, accessed  
853 22-April-2019].
- 854 [79] A. Clark, et al., PIL: Python imaging library, 2010–2019. [Online; accessed 10-April-  
855 2019].
- 856 [80] X. Cao, W. Ren, W. Zuo, X. Guo, H. Foroosh, Scene text deblurring using text-  
857 specific multiscale dictionaries, IEEE Transactions on Image Processing 24 (2015)  
858 1302–1314.
- 859 **Andrés Mafla** obtained a M.Sc. degree in Computer Engineering from Univer-  
860 sitat Autònoma de Barcelona in 2018. Currently, he is a PhD student at the  
861 Computer Vision Center under the supervision of Dr. Dimosthenis Karatzas. His  
862 research interests includes text detection and recognition, scene text image re-  
863 trieval, multi-modal embeddings and scene understanding.
- 864
- 865 **Rubèn Tito** received his Computer Engineering M.Sc. in 2018 from the Universi-  
866 tat Autònoma de Barcelona, and now he is doing his PhD under the supervision of  
867 Dr. Marçal Rossinyol and Dr. Ernest Valveny at the Computer Vision Center. His  
868 main research interests include text recognition, word spotting and multi-modal  
869 embeddings.
- 870
- 871 **Lluís Gómez** obtained his PhD in 2016 at Universitat Autònoma de Barcelona.  
872 Currently he is a TECNIOspring Research Fellow (H2020 Marie Skłodowska-  
873 Curie actions of the European Union) at the Computer Vision Center, Universitat  
874 Autònoma de Barcelona. His research interests include a variety of different topics  
875 in machine learning and computer vision.
- 876

877 **Marçal Rusiñol** received his PhD in 2009 from the Universitat Autònoma de  
878 Barcelona. He is an associate researcher at the Computer Vision Center. His main  
879 research interests include reading systems, information retrieval and performance  
880 evaluation.

881

882 **Ernest Valveny** received the PhD degree in 1999 from the Universitat Autònoma  
883 de Barcelona, where he is an associate professor and member of the Computer  
884 Vision Center. His research interests are on document analysis and pattern recog-  
885 nition, including robust reading, text recognition and retrieval, document classifi-  
886 cation and graph matching.

887

888 **Dimosthenis Karatzas** received his PhD in 2003 from the University of Liver-  
889 pool. He is associate professor at the Universitat Autònoma de Barcelona and  
890 associate director of the Computer Vision Centre where he leads the vision and  
891 language research line. His research interests include reading systems, multi-modal  
892 embeddings, and image captioning.