

Deep Convolutional Image Retrieval: A General Framework

Maria Tzelepi, Anastasios Tefas

Aristotle University of Thessaloniki, Department of Informatics

Abstract

In this paper a Convolutional Neural Network framework for Content Based Image Retrieval is proposed. We employ a deep CNN model to obtain the feature representations from the activations of the deepest layers and we retrain the network in order to produce more efficient image descriptors, relying on the available information. Our method suggests three basic model retraining approaches. That is, the Fully Unsupervised Retraining, if no information except from the dataset itself is available, the Retraining with Relevance Information, if the labels of the dataset are available, and the Relevance Feedback based Retraining, if feedback from users is available. We propose these approaches independently or in a pipeline, where each retraining approach operates as a pretraining step to the subsequent one. We also apply a query expansion method with spatial reranking on top of these approaches in order to boost the retrieval performance. The experimental evaluation on six publicly available image retrieval datasets indicates the effectiveness of the proposed method in learning more efficient representations for the retrieval task, outperforming other CNN-based retrieval techniques, as well as conventional hand-crafted feature-based approaches.

Keywords: Content Based Image Retrieval, Convolutional Neural Networks, Deep Learning, Query Expansion.

Email addresses: mtzelepi@csd.auth.gr (Maria Tzelepi), tefas@aiaa.csd.auth.gr (Anastasios Tefas)

1. Introduction

Information Retrieval (IR) refers to the process of obtaining material (text documents, images, audio etc.) that satisfies a certain information need from large databases [1]. Over the long history of IR, numerous works emerged in the field of text retrieval [2], audio [3], video [4], and image retrieval [5]. Image retrieval is a research area of IR of great scientific interest since 1970s. Earlier studies include manual annotation of images using keywords and searching by text [6]. Due to the difficulties of text-based image retrieval, deriving from the manual annotation of images, that is based on the subjective human perception, and the time and labor requirements of annotation, in 1990s Content Based Image Retrieval (CBIR) has been proposed [7].

The objective of CBIR is to retrieve images that are relevant to a query image from a large collection based on their visual content [8]. A key issue concerning CBIR is to extract meaningful information from raw data in order to eliminate the so-called semantic-gap [9]. The semantic-gap refers to the difference between the low level representations of images and their higher level concepts. While earlier works focus on primitive features that describe the image content such as color, texture, and shape, numerous more recent works have been elaborated on the direction of finding semantically richer image representations. Among the most effective are those that use the Fisher Vector descriptors [10], Vector of Locally Aggregated Descriptors (VLAD) [11] or combine bag-of-words models [12] with local descriptors such as Scale-Invariant Feature Transform (SIFT) [13].

Several recent studies introduce Deep Learning algorithms [14] against the shallow aforementioned approaches to a wide range of computer vision tasks, including image retrieval [15, 16, 17, 18]. The main reasons behind their success are the availability of large annotated datasets, and the GPUs computational power and affordability.

Deep Convolutional Neural Networks (CNN), [19, 20], are considered the more efficient Deep Learning architecture for visual information analysis. CNNs comprise of a number of convolutional and subsampling layers with non-linear neural

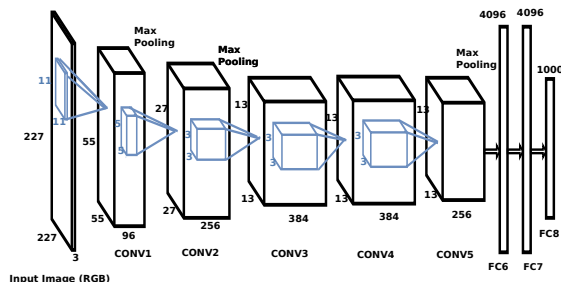


Figure 1: Overview of the CaffeNet Architecture

31 activations, followed by fully connected layers (an overview of the utilized net-
 32 work is provided in Fig. 1). That is, the input image is introduced to the neural
 33 network as a three dimensional tensor with dimensions (i.e., width and height)
 34 equal to the dimensions of the image and depth equal to the number of color
 35 channels (usually three in RGB images). Three dimensional filters are learned
 36 and applied in each layer where convolution is performed and the output is
 37 passed to the neurons of the next layer for non-linear transformation using ap-
 38 propriate activation functions. After multiple convolution layers and subsampling
 39 the structure of the deep architecture changes to fully connected layers and single
 40 dimensional signals. These activations are usually used as deep representations
 41 for classification, clustering or retrieval.

42 Over the last few years, deep CNNs have been established as one of the
 43 most promising avenues of research in the computer vision area due to their
 44 outstanding performance in a series of vision recognition tasks, such as image
 45 classification [21, 22], face recognition [23, 24], digit recognition [25, 26], pose
 46 estimation [27], object and pedestrian detection [28, 29], and action recognition
 47 [30]. It has also been demonstrated that features extracted from the activation
 48 of a CNN trained in a fully supervised fashion on a large, fixed set of object
 49 recognition tasks can be re-purposed to novel generic recognition tasks, [31]. In-
 50 spired by these results, deep CNNs introduced in the vivid research area of CBIR.
 51 The primary approach of applying deep CNNs in the retrieval domain is to ex-
 52 tract the feature representations from a pretrained model by feeding images in

53 the input layer of the model and taking activation values usually drawn from the
54 last layers, while several recent works are directed at utilizing the convolutional
55 layers for the feature extraction. Current research also includes model retraining
56 approaches, which are more relevant to our work, while other studies focus on
57 the combination of the CNN descriptors with conventional descriptors like the
58 VLAD representation. The existing related works are discussed in the following
59 section.

60 Our work investigates model retraining approaches in order to enhance the
61 deep CNN descriptors. We employ a pretrained model to derive feature repre-
62 sentations from the activations of the deepest layers and we retrain the model,
63 exploiting the idea that a deep neural architecture can non-linearly distort the
64 feature space in order to modify the feature representations, with respect to the
65 available information. This information can consist in only the dataset to be
66 searched, the labels of the dataset or of a part of the dataset, and finally infor-
67 mation acquired from users' feedback, that is, relevant or irrelevant images as
68 deemed by multiple users.

69 In this paper we propose a general framework for CNN model retraining in the
70 retrieval domain. The contributions of our study can be summarized as follows:

- 71 • To the best of our knowledge this is the first work that is able to exploit
72 any kind of available information about the retrieval task. The proposed
73 retraining approaches of our method can be categorized as follows:

74 *Fully Unsupervised Retraining (FU)*: if no information is available, except
75 for the dataset itself.

76 *Retraining with Relevance Information (RRI)*: if the labels of the dataset or
77 of a part of the dataset are available.

78 *Relevance Feedback-based Retraining (RF)*: if feedback from users is avail-
79 able.

- 80 • We deploy combinatory schemes, where all the above approaches can be
81 employed in a pipeline. In this fashion each retraining approach operates

82 as a pretraining step to the subsequent one.

- 83 • We suggest a query expansion technique with a spatial verification step
84 applicable to all the above cases.
- 85 • This is the first approach that uses retargeting for the learning phase, instead of triplet loss, allowing for single sample training which is very fast
86 and can be easily parallelized and implemented in a distributed manner.
87

88 In Fig. 2 we schematically describe the proposed framework.

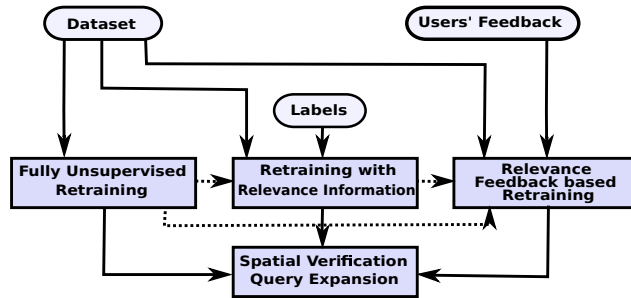


Figure 2: The proposed retraining approaches of our method based on the available information

89 The remainder of the manuscript is structured as follows. Section 2 discusses
90 prior work. The proposed framework is described in detail in Section 3. The pro-
91 posed spatial verification and query expansion technique is presented in Section
92 4. The experiments are provided in Section 5. Finally, conclusions are drawn in
93 Section 6.

94 2. Prior Work

95 In this Section we present previous CNN-based works for image retrieval.
96 Firstly, an evaluation of CNN features in various recognition tasks, including
97 image retrieval that improve the baseline performance using spatial information
98 is presented in [32]. In [33] an image retrieval method, where a CNN pretrained
99 model is retrained on a different dataset with relevant image statistics and classes
100 to the dataset considered at the test time and achieves improved performance, is

101 proposed. From a different viewpoint, in [34, 35], CNN activations at multiple
102 scale levels are combined with the VLAD representation. In [36], a feature aggre-
103 gation pipeline is presented using sum pooling. while in [37] a cross-dimensional
104 weighting and aggregation of deep convolutional neural network layer output is
105 proposed. An approach that produces compact feature vectors derived from the
106 convolutional layer activations that encode several image regions is proposed in
107 [38]. In [39], a three-stream Siamese network is proposed to optimize the weights
108 of the so-called R-MAC representation, proposed in [38], for the retrieval task, us-
109 ing a triplet ranking loss. The public Landmarks dataset, that is also used in [33],
110 is utilized for the model training. In [40] a pipeline that uses the convolutional
111 CNN-features and the bag-of-Words aggregation scheme is proposed. Finally, in
112 [41], the bilinear CNN-based architectures [42] are introduced in the CBIR do-
113 main where a bilinear root pooling is proposed to project the features extracted
114 from the two parallel CNN models into a small dimension and the resulting model
115 is trained on image retrieval datasets using unsupervised training.

116 Subsequently, in [43] an online learning method to learn a similarity func-
117 tion between heterogeneous data modalities by preserving relative similarity con-
118 straints from two directions is proposed. In general, considerable research at-
119 tention has been focused over the past few years on the cross-modal retrieval
120 [44, 45], while another research direction in the retrieval domain, which has at-
121 tracted intensive attention, concerns deep hashing-based techniques [46, 45, 47].
122 Under the hashing view, where the goal is to map the data points into a Ham-
123 ming space of binary codes preserving the similarity in the original space, in
124 [48] a novel unsupervised hashing approach is proposed by integrating feature
125 aggregating and hash function learning into a joint optimization framework. In
126 [45] an end-to-end deep learning framework which can perform feature learning
127 and hash-code learning simultaneously is proposed. Finally, in [46] a two stage
128 hashing framework for cross-modal retrieval tasks which can work in multiple
129 settings like single label, multi-label, and both paired and unpaired scenario, while
130 preserving the structure and semantic relationships that exists within the data is
131 proposed. We should note that the proposed approach can be combined with

132 deep hashing methods to increase the retrieval performance even more, which
133 constitutes a main direction of our future work.

134 A deep CNN is retrained with similarity learning objective function, consider-
135 ing triplets of relevant and irrelevant instances obtained from the fully connected
136 layers of the pretrained model, in [49]. A related approach has also been proposed
137 in the face recognition task which, using a triplet-based loss function, achieves
138 state-of-the-art performance, [50], while a relevant idea recently successfully in-
139 troduced in the cross-modal retrieval domain [51]. These approaches are using
140 triplet sample learning which is difficult to be implemented in large scale, and
141 usually active learning is used in order to select meaningful triplets that can in-
142 deed contribute to learning [50]. In our approach we extend these methodologies
143 by considering multiple relevant and multiple irrelevant samples in the training
144 procedure for each training sample. Additionally, we boost the training speed
145 by defining representation targets for the training samples and regression on the
146 hidden layers, instead of defining more complex loss functions that need three
147 samples for each training step. That is, our approach uses single sample training
148 allowing for very fast and distributed learning. Finally, the proposed method
149 is also able to exploit the geometric structure of the data using unsupervised
150 learning as well as to exploit the user's feedback using relevance feedback.

151 **3. Proposed Method**

152 In this paper we propose a CNN model retraining framework for CBIR, capable
153 of exploiting any kind of available information. The core idea is to utilize the
154 ability of a deep CNN to modify its internal structure, in order to produce better
155 image representations for the retrieval task.

156 We utilize the BVLC Reference CaffeNet model¹, which is an implementation
157 of the AlexNet model trained on the ImageNet Large Scale Visual Recognition
158 Challenge (ILSVRC) 2012 to classify 1.3 million images to 1,000 ImageNet classes,

¹https://github.com/BVLC/caffe/tree/master/models/bvlc_reference_caffenet

159 [21]. The model consists of eight trained neural network layers; the first five
160 are convolutional and the remaining three are fully connected. Max-pooling
161 layers follow the first, second and fifth convolutional layers, while the ReLU non-
162 linearity ($f(x) = \max(0, x)$) is applied to every convolutional and fully connected
163 layer, except the last fully connected layer (denoted as FC8). The output of the
164 FC8 layer is a distribution over 1,000 ImageNet classes. The softmax loss is used
165 during the training. An overview of the CaffeNet architecture is provided in Fig.
166 1.

167 We employ the CaffeNet model to directly extract feature representations from
168 a certain hidden layer. Since the representations obtained from a CNN model for
169 a set of input images are adjustable by modifying the weights of the model, we
170 retrain the parameters of the layer of interest relying on the available information.
171 To this aim, we adapt the pretrained model by removing the layers following the
172 layer utilized for the feature extraction, we build the target representations for
173 each image, and subsequently we retrain the neural network.

174 Based on the available information our method suggests three basic retraining
175 approaches: The FU retraining, if no information is available, the RRI, in the case
176 that the labels of the dataset are available, and the RF, if feedback from users is
177 available. Each of them can be applied independently or in a pipeline, where each
178 approach operates as a pretraining step to the following retraining process. The
179 three basic proposed retraining approaches are presented in detail in the following
180 subsection.

181 *3.1. Model Retraining Approaches*

182 *3.1.1. Fully Unsupervised Retraining*

183 In the FU approach, we aim to amplify the primary retrieval presumption that
184 the relevant images to a certain query are meant to be closer to the query in
185 the feature space. The rationale behind this approach is rooted to the cluster
186 hypothesis which states that documents in the same cluster are likely to satisfy
187 the same information need [52]. That is, we retrain the pretrained CNN model
188 on the given dataset, aiming at minimizing the squared distance between each

189 image representation and its n nearest representations. A schematic description
 190 is provided in Fig. 3.

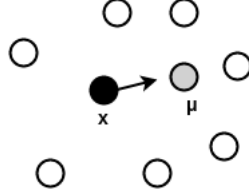


Figure 3: Schematic description of the Fully Unsupervised approach. \circ denote the neighbors of the sample \mathbf{x} , and $\boldsymbol{\mu}$ the mean vector of the nearest neighbors of \mathbf{x}

Let us denote by $\mathcal{I} = \{\mathbf{I}_i, i = 1, \dots, N\}$ the set of N images to be searched, and by $\mathbf{x} = F_L(\mathbf{I})$ the output of the L layer of the pretrained CNN model on an input image \mathbf{I} . Then we denote by $\mathcal{X} = \{\mathbf{x}_i, i = 1, \dots, N\}$ the set of N feature representations emerged in the L layer. We compute the mean vector of the n nearest representations to \mathbf{x}_i and we denote it by $\boldsymbol{\mu}_i$. The new target representations for the images of \mathcal{I} can be determined by solving the following optimization problem:

$$\min_{\mathbf{x}_i \in \mathcal{X}} \mathcal{J} = \min_{\mathbf{x}_i \in \mathcal{X}} \sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu}_i\|_2^2, \quad (1)$$

191 We solve the above optimization problem using gradient descent. The first-
 192 order gradient of the objective function \mathcal{J} is given by:

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial \mathbf{x}_i} &= \frac{\partial}{\partial \mathbf{x}_i} \left(\sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu}_i\|_2^2 \right) \\ &= \frac{\partial}{\partial \mathbf{x}_i} ((\mathbf{x}_i - \boldsymbol{\mu}_i)^\top (\mathbf{x}_i - \boldsymbol{\mu}_i)) \\ &= 2(\mathbf{x}_i - \boldsymbol{\mu}_i), \end{aligned} \quad (2)$$

Consequently, the update rule for the n -th iteration for each image can be formulated as:

$$\mathbf{x}_i^{(n+1)} = \mathbf{x}_i^{(n)} - 2\eta(\mathbf{x}_i^{(n)} - \boldsymbol{\mu}_i), \quad \mathbf{x}_i \in \mathcal{X} \quad (3)$$

193 where the parameter $\eta \in [0, 0.5]$ controls the desired distance from the n nearest
 194 representations.

195 Using the above representations as targets in the layer of interest, we formu-
 196 late a regression task for the neural network, which is initialized on the CaffeNet’s
 197 weights and is trained on the utilized dataset, using back-propagation. The Eu-
 198 clidean loss is used during training for the regression task. Thus, the procedure
 199 is integrated by feeding the entire dataset into the input layer of the modified
 200 model and obtaining the new representations.

201 3.1.2. Retraining with Relevance Information

202 *Samples Provided with Relevance Information.* In this approach we propose to
 203 enhance the performance of the deep CNN descriptors exploiting the relevance
 204 information deriving from the available class labels. To achieve this goal, con-
 205 sidering a labeled representation (\mathbf{x}_i, y_i) , where \mathbf{x}_i is the image representation
 206 and y_i is the corresponding image label, we adapt the deepest neural layers of
 207 the CNN model used for the feature extraction, aiming to minimize the squared
 208 distance between \mathbf{x}_i and the m nearest relevant representations, and simultane-
 209 ously to maximize the squared distance between \mathbf{x}_i and the n nearest irrelevant
 210 representations. A schematic description is provided in Fig. 4.

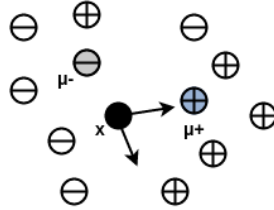


Figure 4: Schematic description of the Supervised approach. \oplus denotes a relevant image to the sample \mathbf{x} , while \ominus denotes an irrelevant one. We indicate the mean vector of relevant images to \mathbf{x} by μ_+ , and the mean vector of irrelevant ones as μ_- .

Let $\mathcal{I} = \{\mathbf{I}_i, i = 1, \dots, N\}$ be a set of N images of the search set provided with relevance information, and $\mathbf{x} = F_L(\mathbf{I})$ the output of the L layer of the pretrained CNN model on an input image \mathbf{I} . Then we denote by $\mathcal{X} = \{\mathbf{x}_i, i = 1, \dots, N\}$ the set of N feature representations emerged in the L layer, by $\mathcal{R}^i = \{\mathbf{r}_k, k = 1, \dots, K^i\}$ the set of K^i relevant representations of the i -th image and by $\mathcal{C}^i = \{\mathbf{c}_l, l = 1, \dots, L^i\}$

the set of L^i irrelevant representations. We compute the mean vector of the m nearest representations of R^i to the certain image representation \mathbf{x}_i , and the mean vector of the n nearest representations of C^i to \mathbf{x}_i , and we denote them by $\boldsymbol{\mu}_+^i$ and $\boldsymbol{\mu}_-^i$, respectively. Then, the new target representations for the images of \mathcal{I} can be determined by solving the following optimization problems:

$$\min_{\mathbf{x}_i \in \mathcal{X}} \mathcal{J}^+ = \min_{\mathbf{x}_i \in \mathcal{X}} \sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu}_+^i\|_2^2, \quad (4)$$

and

$$\max_{\mathbf{x}_i \in \mathcal{X}} \mathcal{J}^- = \max_{\mathbf{x}_i \in \mathcal{X}} \sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu}_-^i\|_2^2. \quad (5)$$

211 We solve the above optimization problems using gradient descent.

The update rules for the n -th iteration can be formulated as:

$$\mathbf{x}_i^{(n+1)} = \mathbf{x}_i^{(n)} - 2\zeta(\mathbf{x}_i^{(n)} - \boldsymbol{\mu}_+^i), \quad \mathbf{x}_i \in \mathcal{X} \quad (6)$$

and

$$\mathbf{x}_i^{(n+1)} = \mathbf{x}_i^{(n)} + 2\beta(\mathbf{x}_i^{(n)} - \boldsymbol{\mu}_-^i), \quad \mathbf{x}_i \in \mathcal{X} \quad (7)$$

Consequently, the combinatory update rule, deriving by adding the equations (6) and (7) can be formulated as:

$$\mathbf{x}_i^{(n+1)} = \mathbf{x}_i^{(n)} - (1 - \beta)(\mathbf{x}_i^{(n)} - \boldsymbol{\mu}_+^i) + \beta(\mathbf{x}_i^{(n)} - \boldsymbol{\mu}_-^i), \quad \mathbf{x}_i \in \mathcal{X} \quad (8)$$

212 where the parameter $\beta = 1 - \zeta, \in [0, 1]$ controls the desired distance both from
 213 relevant and irrelevant representations. Plainly, $\beta = 0$ sets as target representation
 214 for each image the mean vector of its m relevant representations, while as $\beta \rightarrow 1$
 215 the new target representations are more affected by the irrelevant contribution.

216 *Distractors.* In the case where there are images in the dataset that do not be-
 217 long to a certain class and serve as distractors in the retrieval, we can introduce
 218 them to the model retraining procedure. Thus, granted that the distractors are
 219 close to the training samples, that is, their representations are among the n afore-
 220 mentioned irrelevant representations of each image, we concurrently retrain the
 221 pretrained model so that the squared distance between the distractor representa-
 222 tions and each certain image representation be maximized.

223 Denoting by $\mathcal{D} = \{\mathbf{d}_j, j = 1, \dots, P\}$ the set of feature representations of the P
 224 distractors gathered from all the training images, our goal for each corresponding
 225 distractor can be formulated as follows:

$$\max_{\mathbf{d}_j \in \mathcal{D}} \mathcal{J} = \max_{\mathbf{d}_j \in \mathcal{D}} \sum_{j=1}^P \|\mathbf{d}_j - \mathbf{x}_i^j\|_2^2. \quad (9)$$

Consequently, following the gradient, the update rule for the n -th iteration for a
 distractor image can be formulated as:

$$\mathbf{d}_j^{(n+1)} = \mathbf{d}_j^{(n)} + 2\theta(\mathbf{d}_j^{(n)} - \mathbf{x}_i^j), \quad \mathbf{d}_j \in \mathcal{D} \quad (10)$$

226 where the parameter $\theta \in [0, 0.5]$ controls the desired distance from the certain
 227 image representation.

228 Thus, as in the previous approach, using the above target representations we
 229 retrain the neural network on the images provided with relevance information
 230 and on distractors (if any) using back-propagation.

231 3.1.3. Relevance Feedback Based Retraining

232 The idea of this proposed approach is rooted in the relevance feedback phi-
 233 losophy. In general, relevance feedback refers to the ability of users to impart
 234 their judgement regarding the relevance of search results to the system. Then,
 235 the system can use this information to ameliorate its performance [53]. In this
 236 proposed retraining approach we consider information from different users' feed-
 237 back. This information consists of queries and relevant and irrelevant images
 238 to these queries. Then, our goal is to modify the model parameters in order to
 239 bring the relevant images closer to the specific query and move away from it the
 240 irrelevant ones. Towards this end, we retrain the pretrained model by training on
 241 relevant and irrelevant images so that the corresponding relevant representations
 242 come closer in terms of Euclidean distance to the query representation, while the
 243 irrelevant ones move further away. We provide a schematic description in Fig. 5.

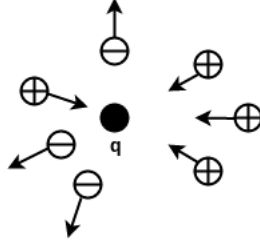


Figure 5: Schematic description of the Relevance Feedback based approach. \oplus denote the relevant images to the query \mathbf{q} , while \ominus denote the irrelevant ones, as they are given by the users

244 Let us denote by $\mathcal{Q} = \{\mathbf{Q}_k, k = 1, \dots, K\}$ a set of queries, $\mathcal{I}_+^k = \{\mathbf{I}_i, i = 1, \dots, N^k\}$
 245 a set of relevant images to a certain query, by $\mathcal{I}_-^k = \{\mathbf{I}_j, j = 1, \dots, M^k\}$ a set of
 246 irrelevant images, by $\mathbf{x} = F_L(\mathbf{I})$ the output of the L layer of the pretrained CNN
 247 model on an input image \mathbf{I} , and by $\mathbf{q} = F_L(\mathbf{Q})$ the output of the L layer on a
 248 query. Then we denote by $\mathcal{X}_+^k = \{\mathbf{x}_i, i = 1, \dots, N^k\}$ the set of feature representations
 249 emerged in L layer of N images that have been qualified as relevant by a user,
 250 and by $\mathcal{X}_-^k = \{\mathbf{x}_j, j = 1, \dots, M^k\}$ the set of M irrelevant feature representations.
 251 The new target representations for the relevant and irrelevant images can be
 252 respectively determined by solving the following optimization problems:

$$\min_{\mathbf{x}_i \in \mathcal{X}_+^k} \mathcal{J}^+ = \min_{\mathbf{x}_i \in \mathcal{X}_+^k} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{q}^k\|_2^2, \quad (11)$$

and

$$\max_{\mathbf{x}_j \in \mathcal{X}_-^k} \mathcal{J}^- = \max_{\mathbf{x}_j \in \mathcal{X}_-^k} \sum_{j=1}^M \|\mathbf{x}_j - \mathbf{q}^k\|_2^2. \quad (12)$$

We solve the above optimization problems using gradient descent. The update rules for the n -th iteration can be formulated as:

$$\mathbf{x}_i^{(n+1)} = \mathbf{x}_i^{(n)} - 2\alpha(\mathbf{x}_i^{(n)} - \mathbf{q}^k), \quad \mathbf{x}_i \in \mathcal{X}_+^k \quad (13)$$

and

$$\mathbf{x}_j^{(n+1)} = \mathbf{x}_j^{(n)} + 2\alpha(\mathbf{x}_j^{(n)} - \mathbf{q}^k), \quad \mathbf{x}_j \in \mathcal{X}_-^k \quad (14)$$

253 where the parameter $\alpha \in [0, 0.5]$ controls the desired distance from the query
 254 representation.

255 Similarly to the other approaches, using the above representations as targets
256 in the layer of interest, we retrain the neural network on the set of relevant
257 and irrelevant images. We note that the above methodology can be implemented
258 in iterative steps as well as in order to improve a certain’s user information
259 need, following the basic relevance feedback concept. More information and
260 experiments are provided in Section 5.5.

261 *3.2. Layer-wise training*

262 The above approaches can be applied on several hidden layers. As mentioned
263 before, several works utilize the fully connected layers [33, 32, 49, 34], since these
264 layers are meant to capture high-level semantic information, while there are also
265 works that utilize the convolutional layers exploiting the spatial information of
266 these layers, using either sum-pooling techniques [36, 37] or max-pooling [38]. In
267 our experiments we apply them on the 7th fully connected layer (FC7), and on
268 the 6th fully connected layer (FC6). The dimension of both the FC6 and FC7
269 layers is 4096 features. Firstly, in the case of the FC7 layer, we employ the
270 CaffeNet model and we adapt it by discarding the FC8 layer and by replacing the
271 ReLU7 layer (that is the ReLU layer following the FC7 layer) with a PReLU layer,
272 [54], which is initialized randomly, and then we retrain it using the appropriate
273 target representations according to the retraining strategy for the CaffeNet’s FC7
274 features. We note that we consider the responses after the ReLU layer. Since
275 the first layers of CaffeNet trained on ImageNet learned more generic feature
276 representations, all the convolutional layers remain unchanged, and we slightly
277 update the FC6 layer using a small learning rate and the FC7 layer with a bigger
278 learning rate, restricting the training cost. In the case of the FC6 modification,
279 we remove the FC7 and FC8 layers, and we replace the ReLU6 layer with a
280 PReLU layer which is initialized randomly, and then we retrain the FC6 layer
281 using proper target representations from the FC6 activation layer of the CaffeNet
282 model.

283 **4. Spatial Verification and Query Expansion**

284 Query Expansion is a standard, in most cases of negligible cost, technique
 285 for accomplishing better retrieval results [55]. The majority of CBIR methods
 286 include a query expansion step that boosts the retrieval performance. On top
 287 of the aforementioned approaches we also introduce a simple query expansion
 288 method by re-issuing the top ten retrieved corresponding image representations
 289 to the initial query as a new query representation, following the average query
 290 expansion scheme.

Let \mathbf{Q} be a certain query, with a CNN representation \mathbf{q} . We consider the
 top ten retrieved images $\mathbf{R}_i, i = 1, \dots, 10$ of \mathbf{Q} and their corresponding CNN
 representations $\mathbf{x}_i, i = 1, \dots, 10$. Then, the new query representation \mathbf{q}_{new} is as
 follows:

$$\mathbf{q}_{new} = \frac{1}{11}(\mathbf{q} + \sum_{i=1}^{10} \mathbf{x}_i). \quad (15)$$

291 Furthermore, we suggest an additional spatial verification step as follows: We
 292 consider a shortlist of N top initially retrieved images for each query \mathbf{Q} , denoted
 293 as $\mathbf{R}_{i,0}, i = 1, \dots, N$. Each of these images is cropped into nine equal-sized over-
 294 lapping regions, $\mathbf{R}_{i,1}, \dots, \mathbf{R}_{i,9}$. An example of the cropping approach is presented
 295 in Fig. 6. Subsequently, we extract the CNN features of the cropped images
 296 and we perform query to the dataset of $N \times 10$ images formed by both the ini-
 297 tial image and the cropped ones. Then, we rerank the shortlist of the initially
 298 retrieved images based on the similarity of the images of the formed dataset to
 299 the query, and we expand the initial query representation as described above with
 300 respect to the reranked list. That is, we rank the $N \times 10$ representations $\mathbf{x}_{i,l}$,
 301 $i = 1, \dots, N, l = 0, \dots, 9$ of the formed dataset in a list, and we perform query ex-
 302 pansion considering the first ten unique corresponding full image representations
 303 of the aforementioned list, $\mathbf{x}_{i,0}$.



Figure 6: Spatial Verification and Query Expansion - An example of image cropping: The first retrieved image, $R_{1,0}$, for a certain query is cropped into 9 overlapping regions denoted as $R_{1,1}, \dots, R_{1,9}$. The height and the width of each region are equal to the half-height and the half-width respectively of the full image

304 5. Experiments

305 In this section we present the experiments conducted in order to assess the
 306 performance of the proposed method. Firstly, a brief description of the evaluation
 307 metrics and the datasets is provided. Subsequently, we describe the experimental
 308 details of each approach, and we demonstrate the experimental results. Finally,
 309 we present the experiments on the proposed relevance feedback technique for a
 310 certain's user information need in iterative steps.

311 5.1. Evaluation Metrics

Throughout this work we use 4 evaluation metrics: precision, recall, mean Average Precision (mAP), and top-N score. The definitions of the above metrics follow below:

$$Precision = \frac{n. \text{ of Relevant Retrieved Images}}{n. \text{ of Retrieved Images}} \quad (16)$$

$$Recall = \frac{n. \text{ of Relevant Retrieved Images}}{n. \text{ of Relevant Images}} \quad (17)$$

Mean Average Precision is the mean value of the Average Precision (AP) of all the queries. The definition of AP for the i -th query is formulated as follows:

$$AP_i = \frac{1}{Q_i} \sum_{n=1}^N \frac{R_i^n}{n} t_n^i, \quad (18)$$

312 where Q_i is the total number of relevant images for the i -th query, N is the total
 313 number of images of the search set, R_i^n is the number of relevant retrieved images
 314 within the n top results; t_n^i is an indicator function with $t_n^i = 1$ if the n -th retrieved
 315 image is relevant to the i -th query, and $t_n^i = 0$ otherwise.

316 Finally, top- N score refers to the average number of same-object images,
 317 within the top- N ranked images.

318 5.2. Datasets

319 **Inria Holidays** [56]: consists of 991 images divided into 500 classes, and
 320 500 discrete queries. Each class in the search set consists of between 1 and 12
 321 images. Some images of the dataset are not in a natural orientation. We note
 322 that we have not proceeded to any preprocessing step of these images, as in other
 323 CNN-based works, e.g. [33, 36]. We measure the retrieval performance in terms
 324 of mAP. Sample images are shown in Fig. 7.



Figure 7: Sample images of the Inria Holidays dataset

325 **Paris 6k** [57]: consists of 6,392 images (20 of the 6,412 provided images are
 326 corrupted) collected from Flickr by searching for particular Paris landmarks and
 327 provides 55 queries. Following the standard evaluation protocol we measure the
 328 retrieval performance in mAP. Like in most CNN-based works [32, 33, 49, 35, 36]

329 we use the full queries for the retrieval. The query images are not considered in
330 the search set in the retrieval procedure, and neither used in the phase of model
331 retraining. We show some example images in Fig. 8.



Figure 8: Sample images of the Paris 6k dataset

332 **UKBench** [58]: contains 10,200 images of objects divided into 2,550 classes.
333 Each class consists of 4 images. All 10,200 images are used as queries. The
334 performance is reported as top-4 score, which is a number between 0 and 4.
335 Samples are provided in Fig. 9.



Figure 9: Sample images of the UKBench dataset

336 **UKBench-2**: since our method performs learning and the UKBench dataset
337 does not provide a discrete set of queries, we hold out one image per class,
338 forming a search set of 7,650 images and a set of 2,550 queries. As in UKBench,
339 we use the top-3 score for the evaluation, which is a number between 0 and 3.

340 **NUS-WIDE** [59]: contains nearly 270,000 images collected from Flickr. It
341 is a multi-label dataset in which each image is annotated with one or multiple
342 concepts from 81 semantic concepts. However, we should note that NUS-WIDE
343 provides links for downloading the images that are not valid, and thus there
344 are differences with the datasets used in previous works. NUS-WIDE dataset is
345 widely used for evaluating hashing techniques for image retrieval, where most of

346 the works (e.g. [60, 61]) utilize the 21 most frequent concepts consisting of at least
347 5,000 images, and the supervised methods use 500 images per concept to form a
348 training set of 10,500 images. In our work, we follow the setting of the 21 most
349 frequent concepts, demanding each image to be associated with only one concept.
350 Thus, we form a database of 40,000 images, with at least 81 images per concept.
351 For each of the 21 concepts we randomly select 100 images, to build the test set
352 of 2,100 queries. We measure the retrieval performance in terms of mAP for the
353 entire used database of 40,000 images. We also report the mAP within the top
354 50 retrieved images. Finally, we use 40,000 additional images that do not belong
355 to any concept and serve as distractors, to test the retrieval performance of our
356 models in the formed database of 80,000 images. Samples are provided in Fig.
357 10.



Figure 10: Sample images of the NUS-WIDE dataset

358 **CIFAR-10** [62]: contains 60,000 images of size 32×32 , divided into 10 classes.
359 Each class contains 6,000 images. Following other works, like [60], we use 50,000
360 images as the dataset to be searched, and we randomly select 1,000 images from
361 the remaining 10,000 images to perform queries. The retrieval performance is
362 measured in terms of mAP, for the entire dataset of 50,000 images. We also
363 report the mAP within the top 50 retrieved images. Sample images are provided
364 in Fig. 11.

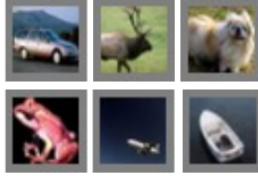


Figure 11: Sample images of the Cifar-10 dataset

365 *5.3. Experimental Setup*

366 The proposed method was implemented using the Caffe Deep Learning frame-
367 work, [63]. In our work we use the adaptive moment estimation algorithm (Adam)
368 [64], instead of the simple gradient descent for the network optimization since it
369 is more stable, with the default parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e - 08$,
370 while the learning rate is set to $1e - 05$. The batch size is set to 64, and the models
371 are trained for 50 epochs. The models are trained on an NVIDIA GeForce GTX
372 1080 with 8GB of GPU memory. All results obtained using Euclidean distance.
373 In the following we present the selected parameters for each of the proposed
374 approaches.

375 *5.3.1. Fully Unsupervised Retraining*

376 In this set of experiments, we consider the 2 nearest representations of each
377 image for the model retraining in all the used datasets, except for the Inria
378 Holidays dataset, where we obtain the new target representations with respect to
379 the 1 nearest representation. The parameter η in (3) is set to 0.5.

380 *5.3.2. Retraining with Relevance Information*

381 In the experiments of this approach, since the number of relevant represen-
382 tations varies meaningfully across datasets, we formulate the new target repre-
383 sentations for the model retraining with respect to each relevant and 5 nearest
384 irrelevant images of each image. The parameter β in (8) is set to 0.2. In Paris 6k
385 dataset, we retrain the network considering relevance information for images an-
386 notated either as good or as ok. Furthermore, for the Paris 6k dataset, where we

Table 1: Inria Holidays

	Scheme	Feature Representation	mAP
1	CaffeNet	CaffeNet \Rightarrow FC6	0.6184
2		CaffeNet \Rightarrow FC7	0.6988
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.6608
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.7307
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.6649
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.74
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.7557
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.7556
9	FU+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.7942
10	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.7687
11	FU+RRI+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.8497

387 utilize information deriving from the available distractors, we note that the num-
388 ber of the utilized distractors varies through datasets and employed approaches.
389 Thus, in Paris dataset, we use information obtained from 846 distractor images
390 for the model retraining, in the FC7 approach (6th row of Table 2), while only 38
391 distractors used in the FC6 one (5th row of Table 2). The parameter θ in (10) for
392 the distractors target formulation is set to 0.5.

393 5.3.3. Relevance Feedback Based Retraining

394 In the experiments that conducted to validate the performance of the Rele-
395 vance Feedback based approach, we consider for each of 500 different users 1
396 relevant and 5 irrelevant images for the Inria Holidays dataset, which forms a
397 training set of 3,000 images. In Paris 6k dataset, 40 relevant and 20 irrele-
398 vant images are considered for each of 55 different users, while in UKBench-2
399 dataset we use 1 relevant and 1 irrelevant images for the 2,550 different users.
400 In CIFAR-10 dataset we use 12 relevant and 1 irrelevant images for the 1,000
401 different users, while in NUS-WIDE we use 5 relevant and 1 irrelevant images
402 for the 2,100 different users. The parameter α in (13), (14) is set to 0.5.

403 5.4. Experimental Results

404 The three proposed retraining approaches can be applied on several hidden
405 layers. Several works utilize the fully connected layers [33, 32, 49, 34], while there

Table 2: Paris 6k

	Scheme	Feature Representation	mAP
1	CaffeNet	CaffeNet \Rightarrow FC6	0.4621
2		CaffeNet \Rightarrow FC7	0.5388
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.6855
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.6984
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.9794
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.9808
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.6418
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.6547
9	FU+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.7714

Table 3: UKBench

	Scheme	Feature Representation	Score
1	CaffeNet	CaffeNet \Rightarrow FC6	3.1308
2		CaffeNet \Rightarrow FC7	3.3501
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	3.48
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	3.5559
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	3.9927
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	3.9371

Table 4: UKBench-2

	Scheme	Feature Representation	Score
1	CaffeNet	CaffeNet \Rightarrow FC6	2.2086
2		CaffeNet \Rightarrow FC7	2.3996
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	2.4345
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	2.5878
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	2.6996
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	2.7769
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	2.3400
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	2.5020
9	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	2.8251
10	FU+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	2.6396

Table 5: NUS-WIDE 40k

	Scheme	Feature Representation	mAP	mAP@50
1	CaffeNet	CaffeNet \Rightarrow FC6	0.0962	0.1608
2		CaffeNet \Rightarrow FC7	0.1276	0.1734
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.114	0.247
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.1606	0.326
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.1532	0.2806
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.2242	0.3521
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.1439	0.3537
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.1856	0.4255
9	FU+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.2095	0.4561
10	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.2350	0.3599

Table 6: NUS-WIDE 80k

	Scheme	Feature Representation	mAP	mAP@50
1	CaffeNet	CaffeNet \Rightarrow FC6	0.0539	0.1403
2		CaffeNet \Rightarrow FC7	0.0772	0.2155
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.064	0.1621
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.094	0.2285
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.092	0.2005
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.1395	0.2529
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.1098	0.32
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.1364	0.3777
9	FU+RF	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.1489	0.4017
10	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.1434	0.2627

Table 7: CIFAR-10 - CaffeNet Initialization

	Scheme	Feature Representation	mAP	mAP@50
1	CaffeNet	CaffeNet \Rightarrow FC6	0.2153	0.4613
2		CaffeNet \Rightarrow FC7	0.2533	0.5210
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.2423	0.4707
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.2862	0.5393
5	RRI	CaffeNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.332	0.5212
6		CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.4297	0.5942
7	RF	CaffeNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.2444	0.5676
8		CaffeNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.2766	0.6232
9	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.4585	0.6044

Table 8: CIFAR-10 - KevinNet Initialization

	Scheme	Feature Representation	mAP	mAP@50
1	KevinNet	KevinNet \Rightarrow FC6	0.2922	0.5988
2		KevinNet \Rightarrow FC7	0.6024	0.847
3	FU	CaffeNet \rightarrow FU(FC6; FC6) \Rightarrow FC6	0.3756	0.6897
4		CaffeNet \rightarrow FU(FC6,FC7; FC7) \Rightarrow FC7	0.6379	0.8466
5	RRI	KevinNet \rightarrow RRI(FC6; FC6) \Rightarrow FC6	0.53	0.72
6		KevinNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.6989	0.8519
7	RF	KevinNet \rightarrow RF(FC6; FC6) \Rightarrow FC6	0.3882	0.76
8		KevinNet \rightarrow RF(FC6,FC7; FC7) \Rightarrow FC7	0.6377	0.8542
9	FU+RRI	CaffeNet \rightarrow FU(FC6,FC7; FC7) \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7	0.7285	0.8532

406 are also works that utilize the convolutional layers with pooling techniques to
 407 produce the image descriptors [36, 37, 38]. In our work we use the fully-connected
 408 layers, since these layers are meant to capture high level semantic information.
 409 Thus, experiments conducted for each of the three proposed retraining approaches
 410 on the FC6 and the FC7 layers (we did not considered the FC8 layer, since it is
 411 a distribution over the 1,000 ImageNet class labels). We note that by utilizing
 412 the FC7 layer, we produce richer descriptors which can lead to better retrieval
 413 performance, since the FC7 layer captures higher level concepts as compared
 414 to the FC6 layer, however it comes with additional computational cost. More
 415 information can be found in the section 5.4.1 that discusses the computational
 416 cost. Furthermore, based on the available information, the retraining approaches
 417 can be applied in a pipeline. In this fashion, each retraining approach operates
 418 as a pretraining step to the subsequent one. For example, the Fully Unsupervised
 419 approach can be applied as pretraining step to both the Retraining with Relevance
 420 Information and the Relevance Feedback based approaches, since it requires no
 421 additional information except for the dataset itself. Therefore, we have conducted
 422 indicative experiments building combinatory retraining schemes, investigating the
 423 assumption that the combinatory schemes can improve the single-step training
 424 approaches.

425 In the following we denote by FC6 and FC7 the feature representations ob-
 426 tained from the FC6 and FC7 layer of the CNN model respectively. We also
 427 abbreviate the applied query expansion technique to QE, and the spatial verifica-
 428 tion and query expansion to SVQE. Finally, we denote by $FU(L_1, L_2, \dots; L_T)$ the
 429 fully unsupervised retraining on the layers L_1, L_2, \dots with target representations
 430 obtained from the L_T layer, by $RRI(L_1, L_2, \dots; L_T)$ the retraining with relevance
 431 information on the layers L_1, L_2, \dots with target representations obtained from the
 432 L_T layer, and correspondingly by $RF(L_1, L_2, \dots; L_T)$ the relevance feedback based
 433 retraining. We use consecutive arrows to describe the retraining pipeline of our
 434 approaches, and the implication arrow to show the final feature representation
 435 employed for the retrieval procedure. Thus, $CaffeNet \Rightarrow FC7$ implies that we
 436 obtain the FC7 representations directly from the CaffeNet model and we use them

437 for the retrieval procedure, while $CaffeNet \rightarrow RRI(FC6,FC7; FC7) \Rightarrow FC7$ de-
438 notes that we formulate the target representations using the features emerged in
439 the FC7 CaffeNet layer and we retrain both the FC6 and FC7 layers of the Caf-
440 feNet, then we extract the FC7 representations of the modified model, and we use
441 them for the retrieval.

442 Tables 1 - 8 summarize the experimental results on all the datasets. The best
443 performance is printed in bold. From the provided results several remarks can be
444 drawn. Firstly, we observe that each retraining approach improves the baseline
445 results of CaffeNet in all the used datasets. We also see that the other proposed
446 methodologies applied on the modified via the FU approach model yield better
447 retrieval results, as compared to the CaffeNet’s employment, in any considered
448 case. In some cases this sequential strategy can lead to outstanding performance,
449 as in the UKBench-2 dataset, where the refined with relevance information model
450 on the fully unsupervised model outperforms any other approach. Hence, we
451 mainly suggest the FU retraining as a pretraining step that can be utilized to boost
452 the performance of the other retraining approaches. Additionally, we observe
453 that the modified in RRI fashion descriptors enhance significantly the baseline
454 results of CaffeNet in all the datasets on both the FC6 and FC7 approaches, and
455 in Paris 6k dataset we can accomplish state-of-the-art performance by a single
456 training step. Furthermore it is shown that by applying the proposed approaches
457 in pipelines we can achieve outstanding performance. In the Inria Holidays
458 dataset (Table 1), we notice that the RF approach is more effective than both
459 the FU and RRI. This is reasonable since the training set of the RF approach
460 (consisting of 3,000 images) is considerably larger than the one of the other two
461 approaches (consisting of 991 images). Furthermore, we can observe in the case
462 of the UKBench-2 dataset (Table 3) that the improvement of the RF approach is
463 not as notable as the FU and RRI ones. We attribute this to the comparatively
464 small training set of the RF approach (5,100 against 10,200 images).

465 Regarding the NUS-WIDE dataset, we first examine the impact of the number
466 of training samples to the retrieval performance. That is, we apply our Fully
467 Unsupervised retraining approach using 2,000, 5,000, 13,000 26,000 and 40,000

468 (that is the entire database) training samples. The experimental results are illus-
469 trated in Fig. 12 . As it is shown, the obtained mAP utilizing 13,000 images is
470 0.1606, while it reaches up to 0.1678 utilizing the entire database. Thus, since
471 the number of training samples also comes with computational cost, a good com-
472 promise is to set the number of training samples to 13,000 images. Therefore, in
473 the following, we utilize 13,000 images from the database to train the proposed
474 models, and we use the entire database in the retrieval stage, for the evalua-
475 tion. Furthermore, we test the performance each of the proposed approaches in
476 the extended version of the dataset with 80,000 images, where we use 40,000
477 additional images, that do not belong to any concept and serve as distractors. In
478 Table 5 we illustrate the experimental results on the NUS-WIDE dataset of 40,000
479 images, and in Table 6 we illustrate the experimental results on the NUS-WIDE
480 dataset of 80,000 images. As we can observe in both cases each of the proposed
481 approaches improves notably the baseline results. The RRI approach applied on
482 the FC7 layer achieves the best performance in the single-step retraining, improv-
483 ing the baseline CaffeNet’s results by 10 mAP points. Regarding the RF approach,
484 we see that the improvement is more significant for the top-50 retrieved images,
485 achieving outstanding performance against the other retraining strategies. This
486 is reasonable since we use the top 6 retrieved images for each of the queries to
487 build the dataset for the model retraining, and hence we expect to better improve
488 the top retrieved images. Concerning the combinatory schemes, as noticed in
489 the other datasets, we observe that we can achieve enhanced performance, as
490 compared to the single-step retraining, by applying the retraining approaches in
491 a pipeline.

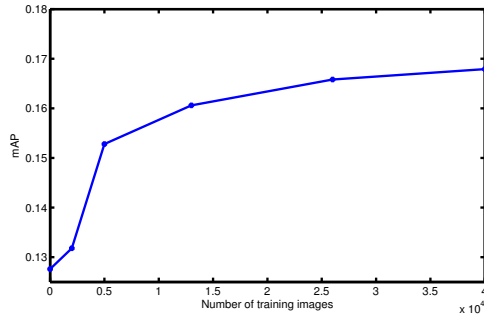


Figure 12: mAP using the FU retraining approach, for various numbers of training samples

492 In the case of the CIFAR-10 dataset apart from the CaffeNet pretrained model,
 493 we also use the KevinNet model [47], which is trained on the CIFAR-10 dataset
 494 for producing binary hash codes for the retrieval task, validating our claim that
 495 the proposed method is applicable to various model architectures, as well as to
 496 finetuned weights for different tasks. We use either the FC7 or the FC6 repre-
 497 sentations, in order to maintain the computational complexity of the proposed
 498 approach, however we could also use the representations produced by the subse-
 499 quent encoding layer, the so called fc8 kevin encode layer. We also tested the
 500 performance of the proposed RRI approach in the aforementioned layer, achieving
 501 a considerable improvement from 0.7907 to 0.8369 in terms of mAP. This also
 502 confirms that the proposed approach can be applied in combination with other
 503 approaches for image retrieval. Furthermore, following this direction, we also
 504 evaluated the hashing codes produced by the RRI optimized fc8 kevin encode
 505 layer, and we report a significant improvement from 0.7863 to 0.8466 in terms of
 506 mAP. As it is shown in Tables 7 and 8 the KevinNet model achieves notably better
 507 baseline results, which is reasonable since it is finetuned on CIFAR-10 dataset.
 508 Furthermore, the observations drawn in the previous datasets, are also confirmed
 509 in CIFAR-10. That is, each of the proposed approaches improve the baseline re-
 510 sults of CaffeNet and KevinNet correspondingly. Additionally, concerning the RF
 511 approach, similarly to the NUS-WIDE dataset, we see that achieves significantly
 512 better results for the top 50 retrieved images, as expected since we use the top
 513 13 retrieved images of each query, and thus the top retrieved images are better

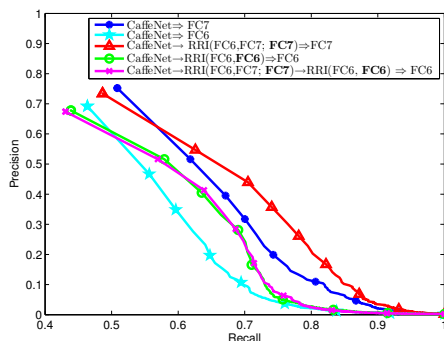


Figure 13: Inria Holidays: Precision-Recall curves of RRI

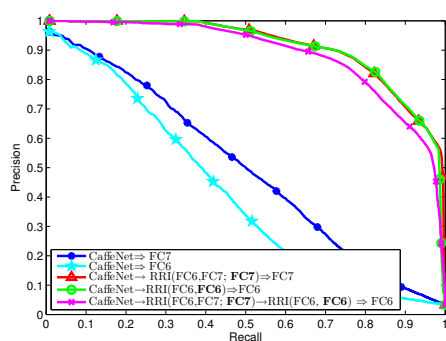


Figure 14: Paris 6k: Precision-Recall curves of RRI

514 improved. Finally, we can observe that the proposed approaches on the KevinNet
 515 initialization, while improve notably the mAP over all the dataset, achieve compar-
 516 atively poorer improvement over the top 50 retrieved images in the case of
 517 the FC7 representations. This is attributed to the fact that the optimized weights
 518 for the binary hashing retrieval task, achieve already enhanced performance, and
 519 thus the proposed method can slightly boost them. On the contrary, we observe
 520 that we can achieve more significant improvement for the 50 retrieved images on
 521 the CaffeNet initialization.

522 In Fig. 13, 14 we provide the Precision-Recall curves of the considered ap-
 523 proaches of the RRI scheme for the Inria Holidays and Paris 6k datasets respec-
 524 tively. In both the datasets, the FC7 modification yields better performance.

525 In Fig. 15, 16 we illustrate the the Precision-Recall curves for the combinatory
 526 schemes on Inria Holidays, and Paris 6k datasets respectively. It is shown that
 527 we can indeed achieve significantly enhanced results by applying our retraining
 528 approaches in a pipeline as compared to the independent ones.

529 In Fig. 17, 18, 19 we provide some examples of the top retrieved images
 530 for certain queries, using the baseline CaffeNet’s features and features obtained
 531 from our retrained models, in Paris 6k, Inria Holidays and UKBench-2 datasets,
 532 respectively.

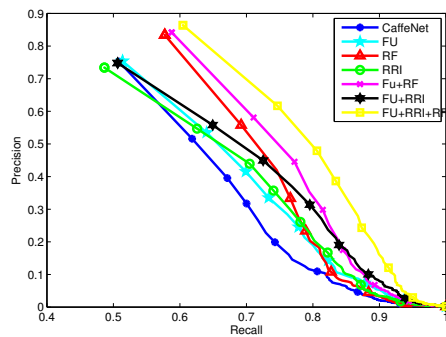


Figure 15: Inria Holidays: Combinatory Schemes

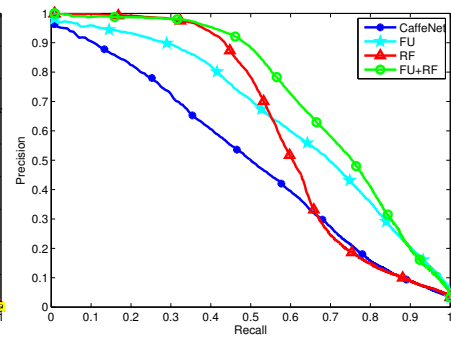


Figure 16: Paris 6k: Combinatory Schemes



Figure 17: Paris: The query image is the first one of the top row and the images that follow in the top row are the first 6 retrieved using the baseline FC7 representation. The top 6 retrieved images using the RRI approach on the FC7 layer are shown in the second row for the same query



Figure 18: Inria Holidays: The query image is the first one of the top row and the images that follow in the top row are the first 5 retrieved using the baseline FC7 representation. The top 5 retrieved images using the RRI approach on the FC7 layer are shown in the second row for the same query



Figure 19: UKBench-2: The query image is the first one of the top row and the images that follow in the top row are the first 3 retrieved using the baseline FC6 representation. The top 3 retrieved images using the RRI approach on the FC6 layer are shown in the second row for the same query

533 5.4.1. Computational Cost

534 The proposed method requires a CNN pretrained model, ideally trained on
 535 the ImageNet dataset composed of 1.2 million images divided into 1,000 classes,
 536 since it produces a rich description of the physical world. Training such a model,
 537 depending on the available GPUs, requires roughly a few days. However, a com-
 538 mon practice in CNN-based works in the retrieval domain is to utilize a pretrained
 539 CNN model, and hence in our work, as previously stated, we utilize the CaffeNet
 540 model. Subsequently, applying each of the proposed methods, requires a certain
 541 training time. Once the models are trained based on the available information,
 542 no additional time is required for the retrieval procedure. That is, the testing
 543 complexity is exactly the same as the baseline models (e.g. CaffeNet, or Kevin-
 544 Net). Regarding the training time, the experiments conducted on an NVIDIA
 545 GeForce GTX 1080 with 8GB of GPU memory, where the average backward pass
 546 time for an input image of the fixed size of 227×227 is 3.32 ms, while the forward
 547 pass takes 2.73 ms, for the model which produces output at the FC7 layer, and
 548 correspondingly, the average backward pass time is 2.53 ms, and the forward one
 549 is 1.97 ms, for the model which produces output at the FC6 layer. Furthermore we
 550 also performed experiments on an NVIDIA GeForce GTX 1060 with 6GB of GPU
 551 memory, as well as on an NVIDIA Quadro K4000 with 3GB of GPU memory, to
 552 measure the training time. The results are illustrated in Table 9.

Table 9: Training time for an input image on various GPUs (FC7 model)

GPU	Backward Pass	Forward Pass
NVIDIA GeForce GTX 1080	3.32 ms	2.73 ms
NVIDIA GeForce GTX 1060	5.86 ms	2.79 ms
NVIDIA Quadro K4000	13.23 ms	9.07 ms

553 In order to improve the deploy speed of the proposed models, we utilize the
 554 NVIDIA TensorRT² tool. TensorRT is a high-performance learning inference
 555 library, which automatically optimizes trained neural networks for run-time per-
 556 formance. Thus, using TensorRT we achieve a significant speed up in both the
 557 proposed model architectures. That is, for the FC7 model the forward pass takes
 558 1.43 ms, while for the FC6 model it takes 1.18 ms.

559

560 *5.4.2. Impact of the probabilistic factors*

561 In this work, we propose a model retraining framework, which is overall able
 562 to exploit any kind of available information. The core idea is that we utilize
 563 a pretrained CNN model, in order to derive the feature representations of a
 564 deep layer, and we retrain the weights of the model, exploiting the idea that
 565 a deep neural architecture can non-linearly distort the feature space in order
 566 to modify the feature representations, with respect to the available information.
 567 Hence, the utilization of fixed weights as the model initialization for the retraining
 568 task, leads to deterministic results, in the retrieval performance. In this section,
 569 we investigate the impact of the probabilistic factors in the performance of the
 570 proposed method. That is, the ordering of input data and the different test images
 571 to perform queries. We choose the CIFAR-10 dataset, to explore the impact of the
 572 aforementioned factors, since the test set of the rest of the datasets is fixed, not
 573 allowing to straightforwardly perform queries with different images. We use as
 574 weights initialization the CaffeNet model. Thus, we repeat each of the FU, RRI,

²<https://developer.nvidia.com/tensorrt>

575 and RF experiments 5 times, using different random shuffling of input images,
 576 and we evaluate the retrieval performance of the corresponding retrained models
 577 on the FC7 layer, using 5-fold cross validation on 5 different test sets of 1,000
 578 queries, randomly selected from the provided test set of 10,000 images. We also
 579 compute the mAP for the 5 five different test sets using the CaffeNet model. The
 580 experimental results for the mean value and the standard deviation of the mAP for
 581 the five runs are illustrated in Table 10. It is evident that the probabilistic factors
 582 do no affect the results significantly, giving quite stable performance among the
 583 runs.

Table 10: CIFAR-10: 5-fold Cross Validation

Retraining method	mAP
CaffeNet	0.2553 ± 0.0051
FU	0.2836 ± 0.0057
RRI	0.4379 ± 0.0087
RF	0.2759 ± 0.0036

584 In Table 11 we provide the experimental evaluation of our spatial verification
 585 and query expansion technique on the best approach of each dataset. In the
 586 UKBench-2 dataset we use 100 queries. From the demonstrated results we can
 587 notice that indeed the query expansion improves the retrieval results, while the
 588 spatial reranking step slightly boosts the initial performance. This is reasonable,
 589 since the spatial verification is more useful on the region based image retrieval,
 590 where we perform queries with a specified region of interest. To this aim, we
 591 employ the cropped-queries versions of Paris 6k dataset, and we apply our RRI
 592 method on the initial CaffeNet’s features. The baseline mAP is 0.5345 for Paris
 593 6k dataset. We note that in this version, the corresponding full images of the
 594 cropped queries are included in the search set. Subsequently, we apply our spatial
 595 verification and query expansion approach on the modified representations. Table
 596 12 illustrates the experimental results.

Table 11: mAP & Score - Spatial Verification & Query Expansion on our best approaches

	Paris 6k	UKBench-2
Best Result	0.9808	2.80
QE	0.9915	2.82
SVQE	0.9916	2.85

Table 12: mAP - Spatial Verification & Query Expansion on Refined with Relevance Information FC7 approach

	Paris 6k Cropped
Initial Result	0.9742
QE	0.9890
SVQE	0.9932

597 Finally, in Table 13 we compare our method against other CNN-based, as well
 598 as hand-crafted feature-based methods, on image retrieval. MAP measures the
 599 retrieval performance in the Inria Holidays and Paris 6k datasets, while top-4
 600 score is used in the case of the UKBench dataset. Since the proposed RF ap-
 601 proach is novel, and the competitive methods do not utilize information derived
 602 from users' feedback, the RF results are reported only in Tables 1-8 and we do not
 603 include them in the comparisons. Methods marked with * use the cropped queries
 604 in Paris 6k dataset. To the best of our knowledge, the proposed approach outper-
 605 forms every other competitive method, in two out of three datasets. We should
 606 note that in Inria Holidays dataset, we can accomplish competitive results only
 607 with the RF approach and the combinatory retraining schemes. This is attributed
 608 to the nature of the dataset. The Inria Holidays dataset is composed of 991 train
 609 images belonging to 500 classes, with each class consisting of between 1 to 12 im-
 610 ages. That is, we have less than 2 images per class, on average. Therefore, since
 611 a key factor of success of the RRI approach is the number of relevant images per
 612 class for the new targets formulation in the retraining process, Inria dataset can
 613 not benefit from it. Furthermore, the number of the train images constitutes in

614 general an important factor in the deep CNN learning. Thus, the train dataset,
 615 consisting of 991 images, in both FU and RRI approaches, is small for achieving
 616 competitive results against other methods. On the contrary, as we can observe
 617 in Table 1, the RF retraining approach, which builds a dataset of 3,000 images,
 618 outperforms the aforementioned proposed approaches, and we can achieve com-
 619 petitive results to the state-of-the-art methods by applying combinatory retraining
 620 strategies.

621 The trained models of the proposed framework are available at: [https://](https://github.com/mtzelepi/framework)
 622 github.com/mtzelepi/framework

Table 13: Comparison against other methods

Method	Inria Holidays	Paris 6k	UKBench
CVLAD* [65]	0.827	-	3.62
VLAD* [11]	0.653	-	-
T-embedding* [66]	0.781	-	-
BOW 200k-D* [67]	0.54	0.46	2.81
Neural Codes [33]	0.793	-	3.56
CNNaug-ss [32]	0.843	0.795	3.644
ReDSL.FC1 [49]	-	0.9474	-
Spoc [36]	0.808	-	3.65
CNN-VLAD [35]	0.84	0.694	-
CRB-CNN-16 [41]	0.854	-	3.56
Deep Image Retrieval [39]	0.907	0.912	-
Deep Image Retrieval & QE [39]	-	0.938	-
R-MAC* [38]	-	0.83	-
R-MAC & QE [38]	-	0.865	-
CroW* [37]	0.849	0.796	-
CroW & QE [37]	-	0.83	-
Ours	0.74	0.9808	3.9927
Ours & QE	-	0.9916	-

623 5.5. Relevance Feedback

624 As mentioned before, the relevance feedback based retraining approach can
625 be materialized in order to improve a certain's user information need in iterative
626 steps [68]. Thus, in each feedback round, the user marks either as relevant or as
627 irrelevant the retrieved images, and the system uses this information to retrain the
628 CNN model according to the methodology described in Section 3.1.3 for multiple
629 users. The relevance feedback procedure is integrated by feeding the images of
630 the given dataset and the query image into the input layer of the modified model
631 and obtaining the new representations. The above process is performed in each
632 relevance feedback round, by initializing the CNN model with the parameters
633 of the previous round and retraining on the new set of relevant and irrelevant
634 images with their corresponding updated targets. Given a new query from the
635 user, the system executes the procedure from the beginning.

636 In order to evaluate the proposed approach, we perform experiments on the
637 Inria Holidays dataset. We obtain the FC7 representations from the CaffeNet
638 model. We consider as search set 991 images and we perform 100 queries from
639 the residue. Each class in the search set consists of between 1 and 12 images. We
640 execute 3 relevance feedback rounds for each query. At each relevance feedback
641 round we use 3 relevant and 3 irrelevant images for the model retraining. The
642 model is trained for 2 epochs at each relevance feedback round. We measure the
643 performance at each feedback round in terms of Precision, and we compute the
644 average precision obtained over all the performed queries. Average precision is
645 measured for the top 3 retrieved images. Experimental results are illustrated in
646 Fig. 20. We can observe that the proposed methodology improves the retrieval
647 performance by the first feedback round.

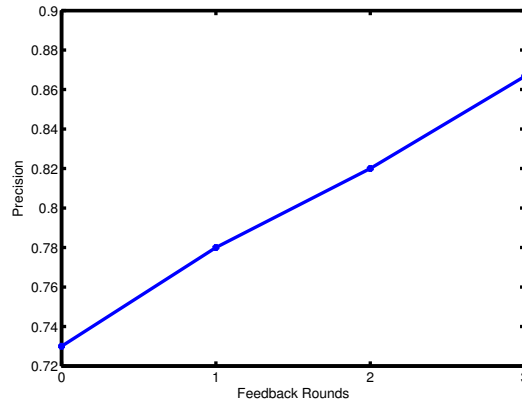


Figure 20: Inria Holidays: Relevance Feedback

648 **6. Conclusions**

649 In this paper we proposed a model retraining framework for enhancing deep
650 CNN representations in the retrieval domain. The proposed method is able to
651 exploit any kind of available information. Thus, if no information is available,
652 the Fully Unsupervised retraining approach is proposed, if the labels are avail-
653 able the Retraining with Relevance Information, and finally if users' feedback
654 is available the Relevance Feedback based retraining is proposed. We utilize a
655 deep CNN model to obtain the feature representations and build the target rep-
656 resentations according to each approach, and then we retrain appropriately the
657 network's weights. We also proposed combinatory retraining strategies, where
658 each of the retraining approaches can be utilized as a pretraining step in order to
659 boost the performance of the following one. A query expansion technique with a
660 spatial verification step applied on top of the best stated approaches provides fur-
661 ther boosting of the retrieval performance. We should note that all the proposed
662 approaches are applicable to any other CNN-based works for image retrieval that
663 utilize a CNN model to directly extract feature representations. Experimental re-
664 sults indicate the effectiveness of our method, performing superior performance
665 over the state of the art approaches, either via a single retraining approach, or by
666 utilizing successive retraining processes.

667 **Acknowledgment**

668 Maria Tzelepi was supported by the General Secretariat for Research and
669 Technology (GSRT) and the Hellenic Foundation for Research and Innovation
670 (HFRI) (PhD Scholarship No. 2826).

671 **References**

- 672 [1] C. D. Manning, P. Raghavan, H. Schutze, Introduction to information re-
673 trieval, Cambridge University Press, 2008.
- 674 [2] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval,
675 Information processing & management 24 (5) (1988) 513–523.
- 676 [3] G. Guo, S. Z. Li, Content-based audio classification and retrieval by support
677 vector machines, IEEE Transactions on Neural Networks 14 (1) (2003) 209–
678 215.
- 679 [4] Y.-G. Jiang, C.-W. Ngo, J. Yang, Towards optimal bag-of-features for object
680 categorization and semantic video retrieval, in: Proceedings of the 6th ACM
681 International Conference on Image and Video Retrieval, ACM, 2007, pp.
682 494–501.
- 683 [5] Y. Rui, T. S. Huang, S.-F. Chang, Image retrieval: Current techniques,
684 promising directions, and open issues, Journal of visual communication and
685 image representation 10 (1) (1999) 39–62.
- 686 [6] N.-S. Chang, K. S. Fu, A relational database system for images, in: Pictorial
687 Information Systems, Springer, 1980, pp. 288–321.
- 688 [7] T. Kato, Database architecture for content-based image retrieval, in:
689 SPIE/IS&T 1992 symposium on electronic imaging: science and technology,
690 International Society for Optics and Photonics, 1992, pp. 112–123.
- 691 [8] R. Datta, J. Li, J. Z. Wang, Content-based image retrieval: approaches and
692 trends of the new age, in: Proceedings of the 7th ACM SIGMM international
693 workshop on Multimedia information retrieval, ACM, 2005, pp. 253–262.

- 694 [9] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, R. Jain, Content-based
695 image retrieval at the end of the early years, *IEEE Transactions on Pattern*
696 *Analysis and Machine Intelligence* 22 (12) (2000) 1349–1380.
- 697 [10] F. Perronnin, Y. Liu, J. Sánchez, H. Poirier, Large-scale image retrieval with
698 compressed fisher vectors, in: *Proceedings of the IEEE Conference on Com-*
699 *puter Vision and Pattern Recognition*, IEEE, 2010, pp. 3384–3391.
- 700 [11] R. Arandjelovic, A. Zisserman, All about vlad, in: *Proceedings of the IEEE*
701 *Conference on Computer Vision and Pattern Recognition*, IEEE, 2013, pp.
702 1578–1585.
- 703 [12] G. Csurka, C. Dance, L. Fan, J. Willamowski, C. Bray, Visual categorization
704 with bags of keypoints, in: *Workshop on statistical learning in computer*
705 *vision, European Conference on Computer Vision (ECCV)*, Vol. 1, Prague,
706 2004, pp. 1–2.
- 707 [13] J. Sivic, A. Zisserman, Video Google: A text retrieval approach to object
708 matching in videos, in: *Proceedings of the International Conference on Com-*
709 *puter Vision*, Vol. 2, 2003, pp. 1470–1477.
- 710 [14] L. Deng, A tutorial survey of architectures, algorithms, and applications for
711 deep learning, *APSIPA Transactions on Signal and Information Processing*
712 3 (2014) e2.
- 713 [15] W. Liu, Z. Wang, X. Liu, N. Zeng, Y. Liu, F. E. Alsaadi, A survey of deep
714 neural network architectures and their applications, *Neurocomputing* 234
715 (2017) 11–26.
- 716 [16] Y. Guo, Y. Liu, A. Oerlemans, S. Lao, S. Wu, M. S. Lew, Deep learning for
717 visual understanding: A review, *Neurocomputing* 187 (2016) 27–48.
- 718 [17] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, A. Y. Ng, Multimodal deep
719 learning, in: *Proceedings of the 28th International Conference on Machine*
720 *Learning (ICML)*, 2011, pp. 689–696.

- 721 [18] Y. Wei, Y. Zhao, C. Lu, S. Wei, L. Liu, Z. Zhu, S. Yan, Cross-modal retrieval
722 with cnn visual features: A new baseline, *IEEE transactions on cybernetics*
723 47 (2) (2017) 449–460.
- 724 [19] B. B. Le Cun, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, L. D.
725 Jackel, Handwritten digit recognition with a back-propagation network, in:
726 *Advances in neural information processing systems 2*, Morgan Kaufmann
727 Publishers Inc., 1990, pp. 396–404.
- 728 [20] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied
729 to document recognition, *Proceedings of the IEEE* 86 (11) (1998) 2278–2324.
- 730 [21] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep
731 convolutional neural networks, in: *Advances in neural information process-*
732 *ing systems*, 2012, pp. 1097–1105.
- 733 [22] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan,
734 V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in: *Proceed-*
735 *ings of the IEEE Conference on Computer Vision and Pattern Recognition*,
736 IEEE, 2015, pp. 1–9.
- 737 [23] Y. Taigman, M. Yang, M. Ranzato, L. Wolf, Deepface: Closing the gap to
738 human-level performance in face verification, in: *Proceedings of the IEEE*
739 *Conference on Computer Vision and Pattern Recognition*, IEEE, 2014, pp.
740 1701–1708.
- 741 [24] D. Wang, J. Yang, J. Deng, Q. Liu, Facehunter: A multi-task convolutional
742 neural network based face detector, *Signal Processing: Image Communica-*
743 *tion* 47 (2016) 476–481.
- 744 [25] D. Ciresan, U. Meier, J. Schmidhuber, Multi-column deep neural networks for
745 image classification, in: *Proceedings of the IEEE Conference on Computer*
746 *Vision and Pattern Recognition*, IEEE, 2012, pp. 3642–3649.
- 747 [26] Y. LeCun, L. Jackel, L. Bottou, C. Cortes, J. S. Denker, H. Drucker, I. Guyon,
748 U. Muller, E. Sackinger, P. Simard, et al., Learning algorithms for classifica-

- 749 tion: A comparison on handwritten digit recognition, *Neural networks: the*
750 *statistical mechanics perspective* 261 (1995) 276.
- 751 [27] A. Toshev, C. Szegedy, Deeppose: Human pose estimation via deep neural
752 networks, in: *Proceedings of the IEEE Conference on Computer Vision and*
753 *Pattern Recognition, IEEE, 2014, pp. 1653–1660.*
- 754 [28] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for
755 accurate object detection and semantic segmentation, in: *Proceedings of the*
756 *IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2014,*
757 *pp. 580–587.*
- 758 [29] D. Tomè, F. Monti, L. Baroffio, L. Bondi, M. Tagliasacchi, S. Tubaro, Deep
759 convolutional neural networks for pedestrian detection, *Signal Processing:*
760 *Image Communication* 47 (2016) 482–489.
- 761 [30] S. Yan, J. S. Smith, B. Zhang, Action recognition from still images based
762 on deep vlad spatial pyramids, *Signal Processing: Image Communication* 54
763 (2017) 118–129.
- 764 [31] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, T. Darrell,
765 Decaf: A deep convolutional activation feature for generic visual recogni-
766 tion., in: *ICML, 2014, pp. 647–655.*
- 767 [32] A. Razavian, H. Azizpour, J. Sullivan, S. Carlsson, Cnn features off-the-
768 shelf: an astounding baseline for recognition, in: *Proceedings of the IEEE*
769 *Conference on Computer Vision and Pattern Recognition Workshops, 2014,*
770 *pp. 806–813.*
- 771 [33] A. Babenko, A. Slesarev, A. Chigorin, V. Lempitsky, Neural codes for image
772 retrieval, in: *European Conference on Computer Vision (ECCV), Springer,*
773 *2014, pp. 584–599.*
- 774 [34] Y. Gong, L. Wang, R. Guo, S. Lazebnik, Multi-scale orderless pooling of deep
775 convolutional activation features, in: *European Conference on Computer*
776 *Vision (ECCV), Springer, 2014, pp. 392–407.*

- 777 [35] J. Ng, F. Yang, L. Davis, Exploiting local features from deep networks for
778 image retrieval, in: Proceedings of the IEEE Conference on Computer Vision
779 and Pattern Recognition Workshops, 2015, pp. 53–61.
- 780 [36] A. Babenko, V. Lempitsky, Aggregating local deep features for image re-
781 trieval, in: Proceedings of the IEEE International Conference on Computer
782 Vision, 2015, pp. 1269–1277.
- 783 [37] Y. Kalantidis, C. Mellina, S. Osindero, Cross-dimensional weighting for ag-
784 gregated deep convolutional features, in: European Conference on Computer
785 Vision (ECCV) Workshops, Springer, 2015, pp. 685–701.
- 786 [38] G. Tolias, R. Sivic, H. Jégou, Particular object retrieval with integral max-
787 pooling of cnn activations, CoRR abs/1511.05879.
- 788 [39] A. Gordo, J. Almazán, J. Revaud, D. Larlus, Deep image retrieval: Learn-
789 ing global representations for image search, in: European Conference on
790 Computer Vision, Springer, 2016, pp. 241–257.
- 791 [40] E. Mohedano, K. McGuinness, N. E. O’Connor, A. Salvador, F. Marques,
792 X. Giro-i Nieto, Bags of local convolutional features for scalable instance
793 search, in: Proceedings of the 2016 ACM on International Conference on
794 Multimedia Retrieval, ACM, 2016, pp. 327–331.
- 795 [41] A. Alzu’bi, A. Amira, N. Ramzan, Content-based image retrieval with compact
796 deep convolutional features, *Neurocomputing* 249 (2017) 95–105.
- 797 [42] T.-Y. Lin, A. RoyChowdhury, S. Maji, Bilinear cnn models for fine-grained
798 visual recognition, in: Proceedings of the IEEE International Conference on
799 Computer Vision, 2015, pp. 1449–1457.
- 800 [43] Y. Wu, S. Wang, Q. Huang, Online asymmetric similarity learning for cross-
801 modal retrieval, in: Proceedings of the IEEE Conference on Computer Vision
802 and Pattern Recognition, 2017, pp. 4269–4278.

- 803 [44] K. Chen, T. Bui, F. Chen, Z. Wang, R. Nevatia, Amc: Attention guided multi-
804 modal correlation learning for image search, arXiv preprint arXiv:1704.00763.
- 805 [45] Q.-Y. Jiang, W.-J. Li, Deep cross-modal hashing.
- 806 [46] D. Mandal, K. N. Chaudhury, S. Biswas, Generalized semantic preserving
807 hashing for n-label cross-modal retrieval, in: Proceedings of the IEEE Con-
808 ference on Computer Vision and Pattern Recognition, 2017, pp. 4076–4084.
- 809 [47] K. Lin, H.-F. Yang, J.-H. Hsiao, C.-S. Chen, Deep learning of binary hash
810 codes for fast image retrieval, in: Proceedings of the IEEE conference on
811 computer vision and pattern recognition workshops, 2015, pp. 27–35.
- 812 [48] T.-T. Do, D.-K. L. Tan, T. T. Pham, N.-M. Cheung, Simultaneous fea-
813 ture aggregating and hashing for large-scale image search, arXiv preprint
814 arXiv:1704.00860.
- 815 [49] J. Wan, D. Wang, S. C. H. Hoi, P. Wu, J. Zhu, Y. Zhang, J. Li, Deep learning
816 for content-based image retrieval: A comprehensive study, in: Proceedings of
817 the ACM International Conference on Multimedia, ACM, 2014, pp. 157–166.
- 818 [50] F. Schroff, D. Kalenichenko, J. Philbin, Facenet: A unified embedding for
819 face recognition and clustering, in: Proceedings of the IEEE Conference on
820 Computer Vision and Pattern Recognition, IEEE, 2015, pp. 815–823.
- 821 [51] L. Wang, Y. Li, S. Lazebnik, Learning deep structure-preserving image-text
822 embeddings, in: Proceedings of the IEEE Conference on Computer Vision
823 and Pattern Recognition, 2016, pp. 5005–5013.
- 824 [52] E. M. Voorhees, The cluster hypothesis revisited, in: Proceedings of the 8th
825 annual international ACM SIGIR conference on Research and development
826 in information retrieval, ACM, 1985, pp. 188–196.
- 827 [53] Y. Rui, T. S. Huang, M. Ortega, S. Mehrotra, Relevance feedback: a power
828 tool for interactive content-based image retrieval, IEEE Transactions on Cir-
829 cuits and Systems for Video Technology 8 (5) (1998) 644–655.

- 830 [54] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing
831 human-level performance on imagenet classification, in: Proceedings of the
832 IEEE International Conference on Computer Vision, 2015, pp. 1026–1034.
- 833 [55] O. Chum, J. Philbin, J. Sivic, M. Isard, A. Zisserman, Total recall: Automatic
834 query expansion with a generative feature model for object retrieval, in: 2007
835 IEEE 11th International Conference on Computer Vision, IEEE, 2007, pp. 1–8.
- 836 [56] H. Jégou, M. Douze, C. Schmid, Hamming embedding and weak geometric
837 consistency for large scale image search, in: A. Z. David Forsyth, Philip Torr
838 (Ed.), European Conference on Computer Vision (ECCV), Vol. I of LNCS,
839 Springer, 2008, pp. 304–317.
- 840 [57] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Lost in quantiza-
841 tion: Improving particular object retrieval in large scale image databases,
842 in: Proceedings of the IEEE Conference on Computer Vision and Pattern
843 Recognition, IEEE, 2008, pp. 1–8.
- 844 [58] D. Nister, H. Stewenius, Scalable recognition with a vocabulary tree, in: Pro-
845 ceedings of the IEEE Conference on Computer Vision and Pattern Recogni-
846 tion, Vol. 2, IEEE, 2006, pp. 2161–2168.
- 847 [59] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, Y. Zheng, Nus-wide: a real-world
848 web image database from national university of singapore, in: Proceedings
849 of the ACM international conference on image and video retrieval, ACM,
850 2009, p. 48.
- 851 [60] H. Lai, Y. Pan, Y. Liu, S. Yan, Simultaneous feature learning and hash
852 coding with deep neural networks, in: Proceedings of the IEEE Conference
853 on Computer Vision and Pattern Recognition, 2015, pp. 3270–3278.
- 854 [61] W.-J. Li, S. Wang, W.-C. Kang, Feature learning based deep supervised
855 hashing with pairwise labels, arXiv preprint arXiv:1511.03855.
- 856 [62] A. Krizhevsky, G. Hinton, Learning multiple layers of features from tiny
857 images.

- 858 [63] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick,
859 S. Guadarrama, T. Darrell, Caffe: Convolutional architecture for fast fea-
860 ture embedding, in: Proceedings of the 22nd ACM international conference
861 on Multimedia, ACM, 2014, pp. 675–678.
- 862 [64] D. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv
863 preprint arXiv:1412.6980.
- 864 [65] W.-L. Zhao, H. Jégou, G. Gravier, Oriented pooling for dense and non-
865 dense rotation-invariant features, in: BMVC-24th British Machine Vision
866 Conference, 2013.
- 867 [66] H. Jégou, A. Zisserman, Triangulation embedding and democratic aggrega-
868 tion for image search, in: Proceedings of the IEEE Conference on Computer
869 Vision and Pattern Recognition, IEEE, 2014, pp. 3310–3317.
- 870 [67] H. Jégou, F. Perronnin, M. Douze, J. Sanchez, P. Perez, C. Schmid, Aggre-
871 gating local image descriptors into compact codes, IEEE Transactions on
872 Pattern Analysis and Machine Intelligence 34 (9) (2012) 1704–1716.
- 873 [68] M. Tzelepi, A. Tefas, Relevance feedback in deep convolutional neural net-
874 works for content based image retrieval, in: Proceedings of the 9th Hellenic
875 Conference on Artificial Intelligence, SETN '16, ACM, 2016, pp. 27:1–27:7.