

DEEP REGIONAL FEATURE POOLING FOR VIDEO MATCHING

Yan Bai^{*,1,2}, Jie Lin^{*,3}, Vijay Chandrasekhar^{*,3,4},
Yihang Lou^{1,2}, Shiqi Wang⁴, Ling-Yu Duan¹, Tiejun Huang¹, Alex Kot⁴

¹Institute of Digital Media, Peking University, Beijing, China

²SECE of Shenzhen Graduate School, Peking University, Shenzhen, China

³Institute for Infocomm Research, A*STAR, Singapore

⁴Rapid-Rich Object Search Laboratory, Nanyang Technological University, Singapore

ABSTRACT

In this work, we study the problem of deep global descriptors for video matching with *regional feature pooling*. We present a theoretical analysis to analyze the joint effect of ROI (Region of Interest) size and pooling moment on video matching performance. To this end, we propose to directly estimate the distribution of matching function scores with the pooling function nested in. Then matching performance can be estimated by the separability of these class-conditional distributions between matching and non-matching pairs. Though the theoretical model is simpler than the video matching and retrieval problem, empirical studies on both pairwise video matching and video retrieval on the challenging MPEG CDVA dataset demonstrate the consistency between our analysis and experimental results.

Index Terms— Convolutional Neural Networks, Pooling, Global Descriptor, Video Matching, Video Retrieval

1. INTRODUCTION

Recent years have witnessed a remarkable growth of interest in video retrieval, which refers to searching for videos representing the same object or scene as the one depicted in a query video. The main challenge is to develop compact and discriminative video feature representations towards highly efficient and effective video matching and retrieval. The motion picture experts group (MPEG) published the standardization of Compact Descriptors for Visual Search (CDVS) [1] in 2015, which came up with a normative bitstream of standardized descriptors for mobile visual search [2] and augmented reality applications [3]. State-of-the-art handcrafted descriptors (VLAD [4], Fisher vectors (FV) [5] and compact FV [6]) built on local invariant SIFTs have been adopted in CDVS as global descriptors. Very recently, MPEG has started a standardization effort titled Compact Descriptors for Video Analysis (CDVA) [7], to extend the CDVS standard to video analysis.

To deal with content redundancy along the temporal dimension, the latest CDVA Experimental Model (CXM) [8] casts video retrieval into keyframe based image retrieval task, in which keyframe-level matching results are combined for video matching. The keyframe-level representation avoids descriptor extraction on dense frames in videos, for reducing computational complexity. Aggregations of local descriptors over video slots have been explored [9, 10]. In this work, we focus on keyframe based approaches.

Though handcrafted descriptors have achieved great success in the CDVS standard [1] and CDVA experimental model, many recent

papers [11–16] have shown the advantage of deep global descriptors for image retrieval, which can be attributed to the remarkable success of Convolutional Neural Networks (CNN) [17, 18]. In particular, state-of-the-art deep global descriptors R-MAC [15] computes the max over a set of Region-of-Interest (ROI) cropped from feature map outputs of intermediate convolutional layers, followed by the average of these regional max-pooled features. Results show that R-MAC offers remarkable improvements over other deep global descriptors like MAC [15] and SPoC [13], while maintaining the same dimensionality.

Previous related work has focused on how to build regional pooled descriptors. Here, we aim to study the relationship between regional feature pooling and video matching. Specifically, we are interested in the key variables, i.e. ROI (Region of Interest) size and pooling moment, which jointly affect matching performance. We make the following contributions,

- We propose directly estimating the distribution of matching function with pooling function nested in. The matching performance can be measured by the separability of these class-conditional distributions between match and non-match sets. The model is capable of predicting *relative* matching performance between different pooling moments, as a function of ROI size.
- Though the model is simple, we verify that the conclusions drawn from the model are largely consistent with quantitative analysis in practical pairwise video matching and video retrieval experiments.
- Experimental results show that deep regional pooled descriptors outperform state-of-the-art deep and handcrafted features on the challenging video matching and retrieval datasets.

Our theoretical analysis is inspired by the work from Boureau et al. [19]. Our major contribution is the analysis of the distribution of matching function scores for matching problem, compared to [6] which focuses on the classification problem. To the best of our knowledge, our analysis is the first attempt to connect deep regional feature pooling to image matching and retrieval performance.

The rest of the paper is organized as follows. In Section 2, we briefly introduce the keyframe based video matching framework, with deep regional pooled global descriptors. Section 3 presents the theoretical analysis of the matching framework. We systematically evaluate the theoretical analysis in Section 4, and summarize the paper in Section 5.

* Yan Bai, Jie Lin, Vijay Chandrasekhar contributed equally.

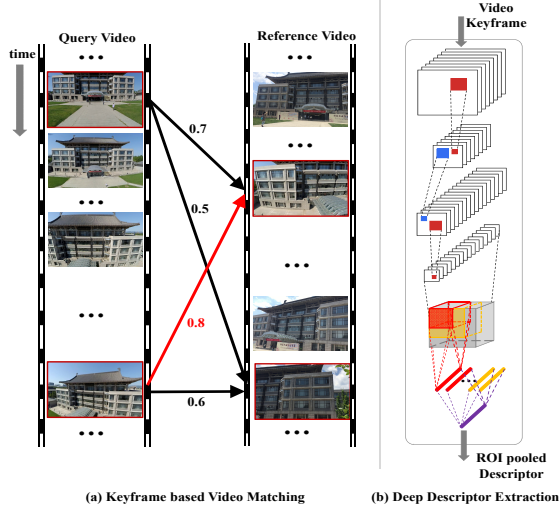


Fig. 1. (a) Keyframe based video matching. (b) Deep regional feature pooling over activation feature maps extracted from CNN architecture.

2. VIDEO MATCHING AND RETRIEVAL

For video matching and retrieval, feature extraction and matching/retrieval processes are included. In this section, we introduce the video matching pipeline and deep regional feature pooling for each keyframe.

2.1. Video Matching and Retrieval Pipeline

Fig. 1 (a) shows the diagram for keyframe based video matching. First, color histogram comparison is applied to identify keyframes in both query and reference videos. Keyframe detection can reduce the temporal redundancy in videos, and reduce the matching complexity between video pairs. Video matching is performed by comparing global descriptors extracted from keyframe pair $\langle \text{query keyframe}, \text{reference keyframe} \rangle$. We denote query video $\mathbf{X} = \{X_1, \dots, X_{N_x}\}$ and reference video $\mathbf{Y} = \{Y_1, \dots, Y_{N_y}\}$, where X and Y denote keyframes. N_x and N_y are the numbers of detected keyframes in query and reference videos, respectively. Each keyframe in \mathbf{X} is compared to all keyframes in \mathbf{Y} . The video-level similarity $K(\mathbf{X}, \mathbf{Y})$ is defined as the largest matching score among all keyframe-level similarity scores.

$$K(\mathbf{X}, \mathbf{Y}) = \max_{X \in \mathbf{X}, Y \in \mathbf{Y}} k(f(X), f(Y)), \quad (1)$$

where $f(X)$ denotes a global descriptor extracted from keyframe X , and $k(\cdot, \cdot)$ represents a matching function.

For video retrieval, the video-level similarity between query and each candidate database video is obtained following the same principle as pairwise video matching. The top ranked candidate database videos are returned for each query video.

2.2. Deep Regional Feature Pooling

Fig. 1 (b) introduces the global descriptor extraction pipeline for each keyframe. Considering an image as input to CNN, we describe it with the feature maps extracted from intermediate convolutional layer, denoted as $\mathbf{X} = \{x_1, \dots, x_C\}$, where x_c represents feature

map of size $W \times H$, and C is the number of channels. ROIs can be sampled from feature maps with varied sizes and strides. For the c^{th} channel, regional feature pooling is first computed by α_s -norm pooling over ROIs sampled from the c^{th} feature map, followed by α_n -norm pooling,

$$f(x_c) = f_{\alpha_n}^{N_{ROI}}(f_{\alpha_s}^{S_{ROI}}(x_c)), \quad (2)$$

where S_{ROI} denotes ROI of size $w \times h$, with $w < W$ and $h < H$. N_{ROI} represents the number of sampled ROIs. $\alpha_s \in \{1, +\infty\}$ and $\alpha_n \in \{1, +\infty\}$ denote pooling moments. For instance, $\alpha = 1$ represents average pooling [13], while $\alpha \rightarrow +\infty$ denotes max pooling [15],

$$f_{\alpha}^N(\hat{x}) = \left(\frac{1}{N} \sum_{i=1}^N (\hat{x}_i)^{\alpha} \right)^{\frac{1}{\alpha}}. \quad (3)$$

The C -dimensional global descriptor $f(\mathbf{X})$ is formed by concatenating $\{f(x_1), \dots, f(x_C)\}$ for all channels.

Following [13], we consider a simple match kernel to compute the similarity between keyframes X and Y with their descriptors $f(X)$ and $f(Y)$,

$$k(f(X), f(Y)) = \beta(X)\beta(Y) \sum_{c=1}^C k(f(x_c), f(y_c)), \quad (4)$$

where $k(f(x_c), f(y_c)) = \langle f(x_c), f(y_c) \rangle$ is the scalar product of pooled descriptors, $\beta(\cdot)$ is a normalization term computed as $\beta(X) = (\sum_{c=1}^C k(f(x_c), f(x_c)))^{-\frac{1}{2}}$. Eq. 4 refers to cosine similarity by accumulating the scalar product of normalized pooled features for each channel. Deep descriptors can be further improved by post-processing techniques such as PCA whitening [13, 15].

3. THEORETICAL ANALYSIS

In this section, we theoretically analyze the effect of ROI pooling on keyframe-level matching performance. More specifically, the matching function Eq. 4 is composed of ROI pooling function Eq. 2, which depends on variables S_{ROI} , N_{ROI} and nested pooling moments (α_s, α_n) . We concentrate our analysis on these variables.

Signal-to-Noise Ratio (SNR) measurement. Inspired by Boureau et al. [19], keyframe-level matching can be regarded as a two-class classification problem (matching and non-matching pairs). As such, we first estimate the distribution (e.g. Binomial distribution with mean μ and variance σ^2) of matching function Eq. 4 performed for each class, respectively. Then, we adopt Signal-to-Noise Ratio (SNR) to measure the separability of these two class-conditional distributions,

$$SNR = \frac{|\dot{\mu}_{k(\cdot, \cdot)} - \ddot{\mu}_{k(\cdot, \cdot)}|}{\dot{\sigma}_{k(\cdot, \cdot)} + \ddot{\sigma}_{k(\cdot, \cdot)}}, \quad (5)$$

where $\dot{\mu}_{k(\cdot, \cdot)}$ ($\ddot{\mu}_{k(\cdot, \cdot)}$) and $\dot{\sigma}_{k(\cdot, \cdot)}$ ($\ddot{\sigma}_{k(\cdot, \cdot)}$) denote mean and standard deviation of Eq. 4 for the match (non-match) set, respectively. The larger the SNR_c , the better the separability between the two classes, implying better matching performance.

Matching function (Eq. 4). For simplicity, we analyze Eq. 4 on per-channel basis, i.e. equal contribution for all channels. Thus, the subscript c is omitted in the following section. Furthermore, by assuming that $f(x)$ are independent identical distribution (i.i.d) random variables, the distribution of matching function, i.e. mean

Table 1. Parameters S_{ROI}, N_{ROI} for ROI pooling function with feature map of size $W \times H = 20 \times 15$.

$S_{ROI} = w \times h (w=h)$	N_{ROI}	$S_{ROI} * N_{ROI}$
15 x 15	3	675
10 x 10	6	600
7 x 7	12	588
3 x 3	65	585

$\mu_{k(\cdot, \cdot)}$ and variance $\sigma_{k(\cdot, \cdot)}^2$, can be computed as follows,

$$\begin{aligned} \mu_{k(\cdot, \cdot)} &= \mu_{f(\mathbf{x})}^2, \\ \sigma_{k(\cdot, \cdot)}^2 &= \sigma_{f(\mathbf{x})}^2(\sigma_{f(\mathbf{x})}^2 + 2 \cdot \mu_{f(\mathbf{x})}^2), \end{aligned} \quad (6)$$

where $\mu_{f(\mathbf{x})}$ and $\sigma_{f(\mathbf{x})}^2$ represent the mean and variance derived from the distribution of ROI pooling function in Eq. 2.

ROI pooling function (Eq. 2). Next, we introduce how to estimate the distribution of ROI pooling function Eq. 2. For simplicity, we consider binary feature maps. We assume all binary spatial bins of \mathbf{x} follow i.i.d. Bernoulli distribution $x_i \sim \text{Bern}(p)$, where p denotes the activation probability of variable x_i . In particular, p is specified as \hat{p} and \tilde{p} for match set and non-match set, respectively.

First, we consider the distribution estimation of Eq.3 without nested pooling moments. For average pooling ($\alpha=1$), the average of N i.i.d Bernoulli variables follows a scaled-down version of Binomial distribution with mean p and variance $\frac{p(1-p)}{N}$. For max pooling ($\alpha \rightarrow +\infty$), the maximum of N i.i.d Bernoulli variables still follows a Bernoulli distribution with mean $(1 - (1-p)^N)$ and variance $(1 - (1-p)^N)(1-p)^N$.

Then, we can derive the distribution of nested pooling moments in a similar way. Specifically, we are interested in the distribution of Eq. 2 with $(\alpha_s, \alpha_n) \in \{(Max - Max), (Avg - Avg), (Max - Avg)\}$. Max-Max follows Bernoulli distribution with

$$\begin{aligned} \mu_{f(\mathbf{x})} &= e(e(p, S_{ROI}), N_{ROI}) \\ \sigma_{f(\mathbf{x})}^2 &= \mu_{f(\mathbf{x})}(1 - \mu_{f(\mathbf{x})}), \end{aligned} \quad (7)$$

Avg-Avg follows Binomial distribution with

$$\begin{aligned} \mu_{f(\mathbf{x})} &= p \\ \sigma_{f(\mathbf{x})}^2 &= \frac{\mu_{f(\mathbf{x})}(1 - \mu_{f(\mathbf{x})})}{S_{ROI} * N_{ROI}}, \end{aligned} \quad (8)$$

Max-Avg follows Binomial distribution with

$$\begin{aligned} \mu_{f(\mathbf{x})} &= e(p, S_{ROI}) \\ \sigma_{f(\mathbf{x})}^2 &= \frac{\mu_{f(\mathbf{x})}(1 - \mu_{f(\mathbf{x})})}{N_{ROI}}, \end{aligned} \quad (9)$$

where function $e(\cdot, \cdot)$ is defined as $e(p, m) = 1 - (1-p)^m$.

Visualization and Conclusion. In summary, we first estimate mean $\mu_{f(\mathbf{x})}$ and variance $\sigma_{f(\mathbf{x})}^2$ from the distribution of ROI pooling function by Eq. 7,8,9, then we obtain mean $\mu_{k(\cdot, \cdot)}$ and variance $\sigma_{k(\cdot, \cdot)}^2$ from the distribution of matching function by Eq. 6. Finally, we compute SNR values following Eq. 5. Note that the former two steps are performed independently for matching and non-matching set.

To visualize SNR values with varied S_{ROI}, N_{ROI} , we simplify our analysis by constraining that $S_{ROI} * N_{ROI} \sim 2 * W * H$ (each spatial bin is roughly scanned twice), in which $S_{ROI} = w \times h$ is

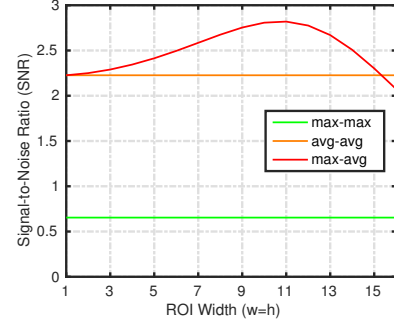


Fig. 2. Video matching performance prediction by measuring the Signal-to-Noise Ratio (SNR) as a function of ROI size, for Max-Max, Avg-Avg, and Max-Avg.

with $w = h$ (i.e. square size ROI), following R-MAC [15]. Table 1 presents examples for S_{ROI}, N_{ROI} with feature map of size $W \times H = 20 \times 15$. Without loss of generality, we empirically set $\hat{p} = 0.02$ and $\tilde{p} = 0.003$, which are statistically computed from match/non-match sets. E.g. a neuron/feature is binarized to 1 if and only if it is activated on both images in a pair, otherwise, 0¹. Correspondingly, we plot SNR as a function of ROI size ($w = h$) for Max-Max, Avg-Avg and Max-Avg, as shown in Fig. 2, which leads us to the main conclusions:

- The overall performance ordering of nested pooling moments is Max-Max is worse than Avg-Avg, while Max-Avg is the best.
- Performance of Max-Max and Avg-Avg is independent on ROI sizes, while Max-Avg firstly rises up and drops later as ROI size increases.

Next, we evaluate the quality of our theoretical analysis with both pairwise video matching and video retrieval experiments on the challenging CDVA dataset.

4. EVALUATIONS

4.1. Datasets

We conduct evaluations on the MPEG CDVA datasets, which contain 9974 query and 5127 reference videos². The videos have durations from 1 sec to 1+ min of 25~30 fps. These videos depict 796 items of interest across three different categories, including 489 large **Landmarks** (5224 queries), 71 **Scenes** (915 queries) and 236 **Objects** (3835 queries, e.g. products). To evaluate the performance of large scale video retrieval, we combine the reference videos with a set of user-generated and broadcast videos as distractors, which consist of content unrelated to the items of interest. There are 14537 distractor videos with more than 1000 hours data. To evaluate pairwise video matching, 4693 match video pairs and 46911 non-match video pairs are constructed from query and reference videos.

We report pairwise matching results in terms of True Positive Rate (TPR), at 1% False Positive Rate (FPR). Video retrieval performance is evaluated by mean Average Precision (mAP).

¹While computation of \hat{p} and \tilde{p} (based on joint activation of neurons in matching and non-matching pairs, respectively) violates independence assumptions, the model still suffices to capture key trends observed for different pooling moments and ROI sizes.

²The MPEG CDVA dataset and evaluation framework are available upon request at <http://www.cldatlas.com/cdva/dataset.html>

Table 2. Pairwise video matching with different combinations of ROI size ($w=h$) and pooling moment, on match/non-match video datasets. No PCA whitening is performed.

ROI size	Max-Max	Avg-Avg	Max-Avg
$w = 15$	71.9	76.8	73.7
$w = 7$			80.7
$w = 3$			77.6

Table 3. Small-scale video retrieval comparison (mAP) with different pooling moments. ROI size is fixed ($w=h=10$). No PCA whitening is performed.

Pool Op.	Landmarks	Scenes	Objects
Max-Max	62.8	79.6	70.5
Avg-Avg	65.3	82.3	69.0
Max-Avg	67.9	83.4	73.8

4.2. Implementation Issues

Following the standard practice in CDVA evaluation framework, there are on average 1~2 keyframes detected per second from raw videos (25~30 fps). We evaluate the theoretical analysis with the widely used VGG16 architecture [18], which is pre-trained on ImageNet ILSVRC classification data set. We resize all video keyframes to VGA size (640×480) images as the inputs of CNN, and extract feature maps of size $20 \times 15 \times 512$ from the last pooling layer (i.e. pool5). Then, we build 512-dimensional deep descriptor following Eq.2. For PCA whitening, we randomly sample 40K frames from the distractor videos for learning the PCA projection.

4.3. Evaluation on Pairwise Video Matching

Table 2 reports pairwise matching performance in terms of True Positive Rate with False Positive Rate equals to 1%, for pooling moments Max-Max, Avg-Avg and Max-Avg, with ROI size $w \in \{15, 7, 3\}$. First, as shown in Table 2, matching performance of Max-Max and Avg-Avg does not vary in ROI size. Second, as ROI size reduced from 15 to 3, matching performance for Max-Avg rises first and drops later. Both observations are consistent with the analysis (Fig. 2). One may note that matching experiments differ from the model in that: (1) The former accumulates similarities contributed by all keyframes and all channels, while the latter is on per-keyframe per-channel basis under i.i.d assumptions. (2) The former works on real-valued feature maps, rather than binary feature maps for the latter case.

4.4. Evaluation on Video Retrieval

Pooling moment. We further design video retrieval experiments to evaluate the quality of the model. Table 3 studies the effect of pooling moments when ROI size is fixed, on all categories of CDVA dataset. We observe the performance order of pooling moments is in line with the analysis in Fig. 2, i.e. Avg-Avg is superior to Max-Max, and Max-Avg performs the best.

ROI size. Table 4 explores the effect of ROI size when pooling moment is fixed (Max-Avg). First, similar to the observation in Table 2, retrieval performance (mAP) increases then decreases as ROI size ranging from 15 to 3. Second, the best performing ROI size depends on the type of data category, i.e. $w = 7$ is the best for Landmarks and Scenes, while $w = 10$ for Objects. Also, there is a quick

Table 4. Small-scale video retrieval comparison (mAP) with different ROI size ($w=h$). Pooling moment is fixed (Max-Avg). No PCA whitening is performed.

ROI size	Landmarks	Scenes	Objects
$w = 15$	60.5	78.6	73.1
$w = 10$	67.9	83.4	73.8
$w = 7$	69.2	84.3	71.4
$w = 3$	64.1	81.3	58.6

Table 5. Large-scale video retrieval comparison of combination (average) of multi ROI sizes with state-of-the-art. The former (latter) number in each cell represents performance without (with) PCA whitening.

Method	Landmarks	Scenes	Objects
CXM [8]	61.4	63.0	92.6
MAC [15]	57.8/61.9	77.4/76.2	70.0/71.8
SPoC [13]	64.0/69.1	82.9/84.0	64.8/70.3
CroW [14]	62.3/63.9	79.2/78.4	71.9/72.0
Max-Avg-Multi	69.4/74.6	84.4/87.3	73.8/78.2

drop in mAP on Objects, e.g. 71.4% ($w = 7$) to 58.6% ($w = 3$). This is probably due to the fact that larger ROI size is desirable for small objects to include contextual info, while it is not that vital for large landmarks and scenes.

Combination of multi ROI sizes. As the best performing ROI size changes across categories, it motivates us to combine multi ROI sizes by averaging their pooled descriptors (termed as Max-Avg-Multi), similar to R-MAC [15]. Table 5 presents the comparison of Max-Avg-Multi against state-of-the-art deep and handcrafted descriptors. Handcrafted descriptors are compact Fisher vectors (FV) built upon SIFT for initial search, followed by geometric reranking with compressed local SIFTs, which has been adopted by the ongoing CDVA standard (terms as CXM). For deep descriptors, we report retrieval performance without (the former number) and with (the latter number) PCA whitening.

We have the following observations: (1) Max-Avg-Multi without PCA whitening achieves the best performance shown in Table 4 on all three categories, which verifies the effectiveness of combination. (2) PCA whitening improves the performance of deep descriptors in most cases. (3) Max-Avg-Multi performs consistently better than other deep descriptors. (4) Overall, deep descriptors outperform handcrafted descriptors by a large margin on Landmarks and Scenes, but are worse on Objects, This is reasonable as handcrafted descriptors based on SIFT are more robust to scale and rotation changes of rigid objects in 2D plane, compared to CNN [16].

5. SUMMARY

In this work, we study the problem of deep regional feature pooling for video matching. In particular, we are interested in the joint effect of ROI size and pooling moment on matching performance. We derive a model to measure matching performance as the separability of distributions of matching function respectively estimated from match and non-match sets. Experiments show that the simplified model is well-aligned with empirical results on both pairwise video matching and video retrieval.

6. REFERENCES

- [1] Ling-Yu Duan, Vijay Chandrasekhar, Jie Chen, Jie Lin, Zhe Wang, Tiejun Huang, Bernd Girod, and Wen Gao, "Overview of the MPEG-CDVS standard," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 179–194, 2016.
- [2] Bernd Girod, Vijay Chandrasekhar, David M Chen, Ngai-Man Cheung, Radek Grzeszczuk, Yuriy Reznik, Gabriel Takacs, Sam S Tsai, and Ramakrishna Vedantham, "Mobile visual search," *IEEE Signal Processing Magazine*, vol. 28, no. 4, pp. 61–76, 2011.
- [3] Mina Makar, Vijay Chandrasekhar, S Tsai, David Chen, and Bernd Girod, "Interframe coding of feature descriptors for mobile augmented reality," *IEEE Transactions on Image Processing*, 2014.
- [4] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez, "Aggregating local descriptors into a compact image representation," in *2010 IEEE Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3304–3311.
- [5] Florent Perronnin, Yan Liu, Jorge Sánchez, and Hervé Poirier, "Large-scale image retrieval with compressed fisher vectors," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [6] Jie Lin, Ling-Yu Duan, Yaping Huang, Siwei Luo, Tiejun Huang, and Wen Gao, "Rate-adaptive compact fisher codes for mobile visual search," *IEEE Signal Processing Letters*, vol. 21, no. 2, pp. 195–198, 2014.
- [7] "Call for Proposals for Compact Descriptors for Video Analysis (CDVA)-Search and Retrieval," *ISO/IEC JTC1/SC29/WG11/N15339*, 2015.
- [8] Werner Bailer Massimo Balestri, Mirosław Bober, "Cdva experimentation model (cxm) 0.2," *ISO/IEC JTC1/SC29/WG11/W16274*, 2015.
- [9] André Araujo, Jason Chaves, Roland Angst, and Bernd Girod, "Temporal aggregation for large-scale query-by-image video retrieval," in *IEEE International Conference on Image Processing*, IEEE, 2015, pp. 1519–1522.
- [10] Zhongwen Xu, Yi Yang, and Alex G Hauptmann, "A discriminative cnn video representation for event detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1798–1807.
- [11] Artem Babenko, Anton Slesarev, Alexandr Chigorin, and Victor Lempitsky, "Neural codes for image retrieval," in *European Conference on Computer Vision*, 2014.
- [12] Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson, "From generic to specific deep representations for visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2015.
- [13] Artem Babenko and Victor Lempitsky, "Aggregating local deep features for image retrieval," in *IEEE International Conference on Computer Vision*, 2015.
- [14] Yannis Kalantidis, Clayton Mellina, and Simon Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *arXiv:1512.04065*, 2015.
- [15] Giorgos Tolias, Ronan Sifre, and Hervé Jégou, "Particular object retrieval with integral max-pooling of cnn activations," in *arXiv:1511.05879*, 2015.
- [16] Vijay Chandrasekhar, Jie Lin, and Olivier Morère, "A practical guide to cnns and fisher vectors for image instance retrieval," *Signal Processing*, vol. 128, pp. 426–439, 2016.
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012.
- [18] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *arXiv:1409.1556*, 2014.
- [19] Y-Lan Boureau, Jean Ponce, and Yann LeCun, "A theoretical analysis of feature pooling in vision algorithms," in *International Conference on Machine learning (ICML)*, 2010.