

# Cross-depiction problem: Recognition and synthesis of photographs and artwork

Peter Hall<sup>1</sup> (✉), Hongping Cai<sup>1</sup>, Qi Wu<sup>2</sup>, and Tadeo Corradi<sup>1</sup>

© The Author(s) 2015. This article is published with open access at Springerlink.com

**Abstract** Cross-depiction is the recognition—and synthesis—of objects whether they are photographed, painted, drawn, etc. It is a significant yet under-researched problem. Emulating the remarkable human ability to recognise and depict objects in an astonishingly wide variety of depictive forms is likely to advance both the foundations and the applications of computer vision. In this paper we motivate the cross-depiction problem, explain why it is difficult, and discuss some current approaches. Our main conclusions are (i) appearance-based recognition systems tend to be over-fitted to one depiction, (ii) models that explicitly encode spatial relations between parts are more robust, and (iii) recognition and non-photorealistic synthesis are related tasks.

**Keywords** cross-depiction; classification; synthesis; feature; spatial layout; connectivity; representation

## 1 Introduction

Many years ago, I took my young children to the zoo. I showed them a simple drawing of a giraffe: bright coloured areas, black lines. When the children got to the zoo, they had no problem at all identifying the giraffe, or the camel, the lion, etc. What is more, they could make recognisable depictions of these animals.

The children were exhibiting (at least) two abilities. One is to generalise from a specific instance to a class, and the other is to generalise

from a depiction (in that case, a particular style of artwork) to real life. Children generalise equally well across depictions; they would have recognised photographs of the animals equally well. Humans are able to recognise objects in an astonishing variety of forms. Whether photographed, drawn, painted, carved in wood, people can recognise horses, bicycles, people, etc. Furthermore, the ability to draw and paint—even from memory—is a strong indicator that in humans at least, recognition and synthesis are related.

The ability of humans to recognise regardless of depiction is such an everyday occurrence that it can often pass without being noticed. Yet it is an astonishing ability that cannot be matched by any current algorithm. Even the very best recognition algorithms—including deep learning—fail to cope with the cross-depiction problem. Indeed, all algorithms we have empirically tested exhibit the same general behaviour: all show a significant drop in performance when presented with an inhomogeneous data set, and fall further still when trying to recognise a drawn object after being trained only on photographic examples. Some algorithms are more pronounced than others in this trend—those that explicitly encode spatial relations tend to be more robust.

The inability of all contemporary approaches to cope with the cross-depiction problem is a significant literature gap. Cross-depiction forces one to consider which visual attributes are necessary for recognition, and which are merely sufficient. That is, one may sensibly ask: which properties of an object class are invariant (or close to invariant) given over variations in depictive style? The specific appearance among different depictive styles varies to a much greater degree than that due to lighting changes, but still

<sup>1</sup> Department of Computer Science, University of Bath, UK. E-mail: maspmh@bath.ac.uk (✉), H.Cai@bath.ac.uk, T.M.Corradi@bath.ac.uk.

<sup>2</sup> School of Computer Science, University of Adelaide, Australia. E-mail: qi.wu01@adelaide.edu.au.

Manuscript received: 2015-03-25; accepted: 2015-05-20

people can recognise them. Children's drawings, as in Fig. 1, are both highly abstract and highly variable, yet contain sufficient information for objects to be recognised by humans, but not computers. Equally overlooked is the fact that no computer is yet able to draw as a child.

Learning the specifics of each depiction seems at best unappealing, not least because the gamut of possible depictions is potentially unlimited. Rather, the question is: what abstraction do these classes have in common that allow them to be recognised regardless of depiction? It is an unavoidable question that pushes at the foundations of computer vision.

A machine that is able to recognise regardless of depiction would provide a significant boost to current applications, such as image search and rendering. For example, given a photograph of the Queen of England, a search should return all portraits of her, including postage stamps that capture her likeness in bas-relief. Searching heterogeneous data sets is a real problem for the creative industries, because they archive vast quantities of material in a huge variety of depictions—a problem that requires visual class models to span depictive styles. Non-photorealistic rendering from images and video would be boosted too, not least because highly aesthetic renderings depend critically on the level of abstraction available to algorithms. Picture making is nothing like tracing over photographs: humans draw what they know of an object, not what they see—computers should do like wise.

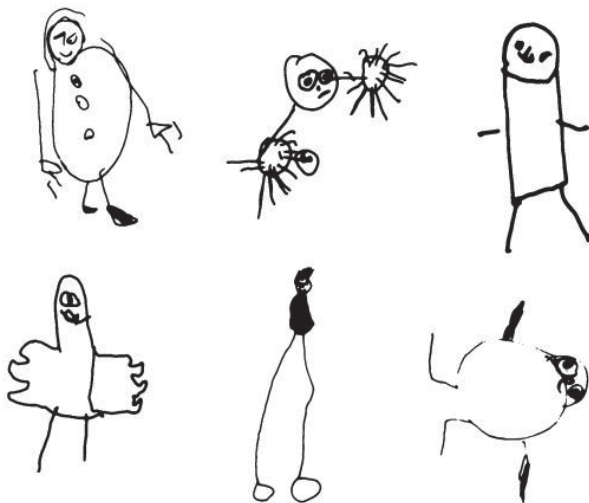


Fig. 1 Children's drawings.

One of our guiding principles has been that the cross-depiction problem acts to unite the synthesis and analysis of images. The rationale is that people find it at best very difficult to draw objects they cannot recognise; more exactly, people tend to draw objects they can see in a manner that is highly influenced by what they know of them. This is most obvious in children, who draw the sky at the top of their pictures and eyes at the top of heads, and often they will draw cars with four wheels, and so on. But it is evident in the artwork of adults too. Indeed, students at western art schools are given extensive life-drawing classes with the exact purpose of teaching them to draw what is seen rather than what is known. For example, early students often draw the hands, feet, and faces in proportion to direct measures rather than as seen when foreshortened: the students' knowledge allows them to compensate for perspective effects.

The key for computational emulation of the human ability is, we argue, representation. It is reasonable to seek a single representation that supports both the recognition and the synthesis of objects. Even so, from an "engineering" point of view, the problems of recognition and synthesis seem sufficiently far apart that different representations are needed. Therefore, we will consider representations that are suitable for each, and then conjecture as to what a single representation might look like.

In summary then, there are two important reasons to study the cross-depiction:

- 1) The "foundational" problem: we are forced to think very carefully about how to model object classes.
- 2) The "practical" consequences: solving the cross-depiction problem will open many robust applications in web search, computer graphics, and other areas.

This paper establishes there is a literature gap; it shows that feature based approaches alone are not sufficient for cross-depiction, and representations that take connectivity and spatial layout into account perform better. It suggests future avenues in terms of object class representation. As a note: in this paper, we use the term photograph as a short hand for "natural image", and the term artwork as all other images.

## 2 Related literature

The computer vision literature distinguishes between classification (does this image contain an object of class  $X$  or not?) and detection (an object of class  $X$  is at this place in this image). Yet lay language makes no sharp distinction; we use the term recognition to mean both classification and detection, which is closer to lay usage.

There is a vast literature in computer vision to address recognition. Yet almost no prior art addresses the cross-depiction, which is surprising given its genuine potential for advancing computer vision both in its foundations and in its applications.

Of the many approaches to visual object classification, the bag-of-words (BoW) family [1–3] is amongst the most widespread. It models visual object classes as histograms of visual words; these words are being clusters in feature space. Although the BoW methods address many difficult issues, they tend to generalise poorly across depictive styles (see Section 3). Alternative low-level features such as edgelets [4, 5] may be considered, or mid-level features such as region shapes [6, 7]. These features offer a little more robustness, but only if the silhouette shape is constrained — and only if the picture offers discernible edges, which is not the case for many artistic pictures (Turner’s paintings, for example).

Deformable models of various types are widely used to model object classes for detection tasks, including several kinds of deformable template models [8, 9] and a variety of part-based models [10–16]. In the constellation models from Ref. [14], parts are constrained to be in a sparse set of locations, and their geometric arrangement is captured by a Gaussian distribution. In contrast, pictorial structure models [12, 13, 15] define a matching problem where parts have an individual match cost in a dense set of locations, and their geometric arrangement is captured by a set of spring connecting pairs of parts. In those methods, the deformable part-based model (DPM) [12] is widely used. It describes an object detection system based on mixtures of multi-scale deformable part models plus a root model. By modelling objects from different views with distinct models, it is able to cope with large variations in pose. None of these directly address the cross-depiction problem.

Shape has also been considered. Leordeanu et al. [17] encode relations between all pairs of edgels of shape to go beyond individual edgels. Similarly, Elidan et al. [18] use pairwise spatial relations between landmark points. Ferrari et al. [19] propose a family of scale invariant local shape features formed by short chains of connected contour segments. Shape skeletons are the dual of shape boundary, and also have been used as a descriptor. For example, Rom and Medioni [20] suggest a hierarchical approach for shape description, combining local and global information to obtain skeleton of shape. Sundar et al. [21] use skeletal graph to represent shape and use graph matching techniques to match and compare skeletons. Shock graph [22] is derived from skeleton models of shapes, and focuses on the properties of the surrounding shape. Shock graphs are obtained as a combination of singularities that arise during the evolution of a grassfire transform on any given shape. In particular, the set of singularities consists of corners, lines, bridges, and other similar features. Shock graphs are then organised into shock trees to provide a rich description of the shape.

Algorithms usually assume that the training and test data are drawn from the same distribution. This assumption may be breached in real-world applications, leading to domain-adaptation methods such as transfer component analysis (TCA) [23], which transfer components from one domain to another. Both sampling geodesic flow (SGF) [24] and geodesic flow kernel (GFK) [4] use intermediate subspaces on the geodesic flow connecting the source and target domain. GFK represents state-of-the-art performance on the standard cross-domain dataset [25]; it has been used to classify photographs acquired under different environmental conditions, at different times, or from different viewpoints.

Cross-depiction problems are comparatively less well explored. Some work is very specific — Crowley and Zisserman take a weakly supervised approach, using a DPM to learn figurative art on Greek vases [26]. Others develop the problem of searching a database of photographs based on a sketch query; edge-based HOG was explored in Ref. [27] and by Li et al. [28]. Others have investigated sketch based retrieval of video [29, 30].

Approaches to the more general cross-depiction

problem are rare. Matching visually similar images has been addressed using self similarity descriptors [31]. It relies on a spatial map built from correlations of small patches; it therefore encodes a spatial distribution, but tends to be limited to small rigid objects. Crowley and Zisserman [32] provide the only example of domain adaptation we know of specifically designed for the cross-depiction problem; they train on photographs and then use midlevel patches to learn spatial consistencies (scale and translation) that allow matching from photographs into artwork. Their method performs well in retrieval tasks for 11 object classes in databases of paintings.

Classification, rather than matching, has also been studied. Shrivista et al. [33] show that an exemplar SVM trained on a huge database is capable of classification of both photographs and artwork. A less computationally intensive approach has been proposed [34] using a hierarchical graph model to obtain a coarse-to-fine arrangement of parts with nodes labelled by qualitative shape [35]. Wu et al. [36] address the cross-depiction problem using a deformable model; they use a fully connected graph with learned weights on nodes (the importance of nodes to discriminative classification), on edges (by analogy, the stiffness of a spring connecting parts), and multiple node labels (to account to different depictions); a method tested on 50 categories. Others use no labels at all, but rely on connection structure alone [37] or distances between low-level parts [17].

Deep learning has recently emerged as a truly significant development in computer vision. It has been successful on conventional databases, and over a wide range of tasks, with recognition rates in excess of 90%. Deep learning has been used for the cross-depiction problem, but its success is less clear cut. Crowley and Zisserman [38] are able to retrieve paintings in 10 classes at a success rate that does not rise above 55%; their classes do not include people. Ginosar et al. [39] use deep learning for detecting people in Picasso paintings, achieving rates of about 10%.

Other than this paper, we know of only two studies assessing the performance of well established methods on the cross-depiction problem. Crowley and Zisserman [32] use a subset of the “Your Paintings” dataset [40], the subset decided by those

that have been tagged with VOC categories [41]. Using 11 classes, and objects that can only scale and translate, they report an overall drop in per class Prec@k (at  $k = 5$ ) from 0.98 when trained and tested on paintings alone, to 0.66 when trained on photographs and tested on paintings. Hu and Collomosse [27] use 33 shape categories in Flickr to compare a range of descriptors SIFT, multi-resolution HOG, Self Similarity, Shape Context, Structure Tensor, and (their contribution) Gradient Field HOG. They test a collection of 8 distance measures, reporting low mean average precision rates in all cases.

Regarding synthesis, non-photorealistic rendering from photographs is germane to our paper. Almost all of the non-photorealistic rendering (NPR) from photographs literature concerns the development of image filtering of one kind or another (see for example Ref. [42] for a review). However, such algorithms fail to emulate the process of human produced arts, which is inevitably about abstraction of some kind, meaning a summary of the object or scene being drawn. Moreover, humans can and do draw (and paint) from memory.

### 3 Representations for recognition

Here we will consider representations for recognition of object classes, regardless of how they are depicted. We describe representations we have used, and benchmark some of them against datasets we have created.

#### 3.1 Feature based representations

As already mentioned in Section 2, bag-of-words (BoW) models for object classes are widespread. BoW models are premised on the assumption that object classes can be distinguished from the relative proportion of discriminative image patches in an unordered collection. Since “words” in the context of images means an image patch, the consequence of the this assumption is that words in patch must exhibit low variation — they must be similar.

Intuitively, this “BoW assumption” is violated when the datasets contain both photographs and artwork; our intuition is confirmed by experiments. In order to see how the local features affect the performance in cross-depiction classification, we

test a range of different features, e.g., SIFT [43], Geometric Blur (GB) [44], Self-similarity Descriptors (SSD) [45], Histogram of Oriented Gradient (HOG) [46], and Edge-based HOG (eHOG) [50].

The BoW we use is the spatial pyramid [2], as it is well known and widely used. Given a set of labelled training images, local descriptors are computed on a regular grid with multiple-sized regions. A vocabulary of words is constructed by vector quantisation of local descriptors with  $k$ -means clustering ( $k = 1000$ ). To construct a visual class model (VCM), each image is partitioned into  $L$  levels of increasingly fine cells ( $L = 2$  in our experiments). A histogram of word occurrences is built for each cell; concatenating these histograms encodes the image with a 5000 dimensional vector. A one-versus-all linear SVM classifier is trained on a  $\chi^2$ -homogeneous kernel map [47] of all training histograms. Given a test image, the local features are extracted in the same way as in the training stage, mapped onto the codebook to build a multi-

resolution histogram, which is then classified with the trained SVM.

We evaluate the algorithms on Photo-Art-50 dataset [36] which contains 50 distinct object classes (see Fig. 2), with between 90 and 138 images for each class. Each class is approximately half photographs and half artwork. All 50 classes appear in Caltech-256; a few also appear in PASCAL VOC Challenge [41] and ETH-Shape dataset [48].

As can be seen in Table 1, none of the BoW methods perform well in recognition over a heterogeneous database as ours. We also used Fisher Vectors (FV) [49], which instead describe the distribution of statistics of local features inside each cluster. Consistent with the observation in Ref. [49], it outperforms BoW-SIFT by 2%–3% in all “train–test” settings. In spite of such an improvement, FV still suffers from significant performance drop in the condition of different training and test depiction domains.

In summary, all methods exhibit comparably



**Fig. 2** Photo-Art-50 dataset [36] containing 50 object categories. Each category is displayed with one art image and one photo image.

**Table 1** Classification using feature based representations. Each row is a train/test pattern: Art, Photo, Mixed. Each column is an algorithm with feature, divided into groups: BoW [31, 43, 44, 46, 49, 50], Fisher Vectors [49]. Domain Adaption using GFK [4] has two variants (PCA and LDA), also Subspace Alignment (SA) [25]. Each cell shows the mean of 5 randomized trials. The standard deviation on any column never rises above 2%. Domain-Adaptation methods were tested only on cross-domain train/text patterns

Model		BoW					FV	GFK_PCA	GFK_LDA	SA
Train	Test	SIFT	GB	SSD	HOG	eHOG	SIFT	SIFT	SIFT	SIFT
P	P	84	77	66	72	70	87	—	—	—
M	P	80	72	58	65	63	84	—	—	—
A	P	64	60	39	42	50	66	48	50	45
A	A	74	72	49	55	60	77	—	—	—
M	A	69	67	45	50	56	73	—	—	—
P	A	44	50	31	29	40	47	31	32	29

high performance with homogeneous data comply with the “low variation” assumption (good for photographs) but show a fall when faced with heterogeneous data (photographs and artwork). The fall is most distinct when BoW and Fischer Vectors are trained on photographs and tested on artwork—suggesting the representation is over-fitted to photographic data. Due to the very different distribution of photo and art domains, it is natural to resort to the domain adaptation techniques. In the following, we will investigate how well the domain adaptation could bridge the gap.

Domain adaptation is a process by which a representation built initially for one domain is allowed to somehow adapt to cover a second. Some may say that photographs and artwork belong to different domains, so that domain adaptation may overcome the problems we see with BoW and Fischer Vectors.

Excellent domain adaptive methods include, but are not limited to Refs. [4, 24, 25, 51, 52]. They show clear benefits for photographs captured under different conditions. We tested some of these (details below) using photographs as a source domain for the initial model, which we adapted to the target domain of artwork. Table 1 shows this case to be the most difficult for BoW and Fischer Vectors. We also tested adaptation in the reverse direction (from art to photographs, still difficult for BoW and FV).

Specifically, we implemented and tested two variants of Geodesic Flow Kernel (GFK) [4]: GFK\_PCA projects original features in both domains (source photograph and target artwork) onto a 49

dimensional subspace via with PCA; GFK\_LDA uses supervised dimensionality reduction via linear discriminant analysis—on the source domain only. Subspace Alignment (SA) [25] project  $\mathcal{S}$  and  $\mathcal{T}$  to respective subspaces. Then, a linear transformation function is learned to align the two domains.

The results for these three methods are shown in Table 1. They suggest that domain adaptation using feature representations are not effective.

### 3.2 Models with spatial and structural information

As Table 1 shows, feature based representations are poorly suited to the task of recognition in the cross-domain problem; even domain adaptation proves ineffective. This section describes representations that take spatial and structural relations into account.

We have used structure alone as a representation [37]. Each class representation was a spatially weighted graph built by hierarchical agglomeration, filtered by Laplacian graph energy [53]. Tests using thirteen different classes in a heterogeneous database showed an accuracy (the diagonal of a confusion matrix) of above 85%. This suggests structural and spatial relations are important to cross-depiction; but the experiments are too limited to be conclusive and later tests on a larger dataset in Ref. [34] yields accuracies of around 20% (see Table 2). This suggests space and structure are important, but are insufficiently rich.

Given that proposition that features should not be limited by the statistics of any one domain (e.g.,

**Table 2** Classification using shape and structure. (a) Single domain task, (b) single cross-depiction task, (c) single to mixture depiction task, and (d) mixture cross-depiction task. The character “p” is “photos”, “a” is “art”, and “m” is “mixture”. Dense SIFT was computed using Ref. [5], and structure only followed Ref. [37]

(a)			(b)				
Case 1: training	5p	5a	Case 2: training	8p	10p	8a	10a
Case 1: testing	15p	15a	Case 2 : testing	15a	15a	15p	15p
Dense SIFT	70%	59%	Dense SIFT	43%	47%	49%	51%
Structure only	16%	19%	Structure only	19%	23%	22%	25%
Proposed method	61%	62%	Proposed method	63%	64%	64%	67%

(c)					(d)		
Case 3: training	3a	5a	3p	5p	Case 4: training	6m	10m
Case 3: testing	30m	30m	30m	30m	Case 4: testing	30m	30m
Dense SIFT	46%	50%	50%	54%	Dense SIFT	60%	61%
Structure only	13%	16%	14%	16%	Structure only	21%	24%
Proposed method	58%	61%	56%	61%	Proposed method	62%	65%

photograph, pencil drawing), we next considered simple shapes as features to label a graph. Using shapes as features was inspired by observing the great artists such as Picasso, who construct recognisable objects from circles, squares, and such like.

We first learnt shapes from image segmentation [54] using a fully unsupervised approach, because we wanted to find out whether simple shapes exist in image segmentation independently of human bias. Our algorithm discovered simple shapes that can be named—circle, square, etc. These are seen in Fig. 3. The same figure shows a scale-based hierarchical decomposition of an image with segments classified using these shapes, plus a “noise” category for segments that did not classify into any shape. A mean graph was used to connect shapes in each layer of the hierarchy, also in Fig. 3. Edges also connect corresponding nodes between layers.

This was tested on a smaller image data base than in Section 3.1, and compared with dense SIFT [5] and structure only [37]. This representation maintains performance across domains—that is, it does not exhibit a fall-off when trained on one domain and tested on another, and all others do so far. Even so, a classification rate hovering around 60% cannot be

regarded as satisfactory: we must turn to stronger models.

### 3.3 DPM, ADPM, and multi-label graph

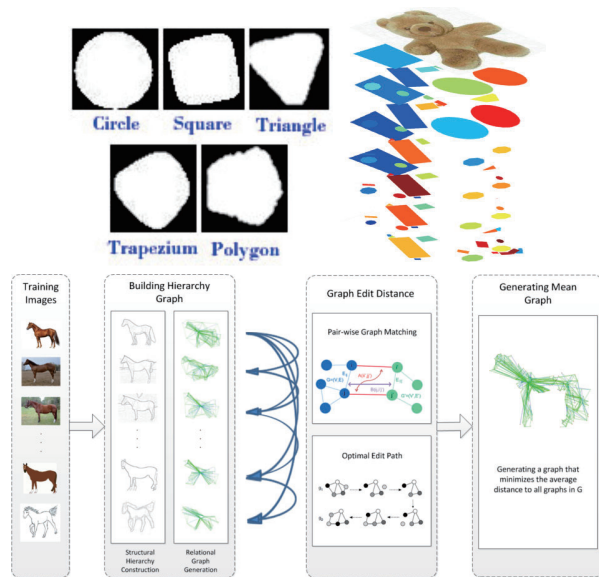
Deformable parts model (DPM) [55] is a well known object representation that takes spatial layout into account. It models an object with a star graph, i.e., a root filter plus a set of parts. Given the location of the root and the relative location of  $n$  parts,  $n = 8$  in our experiments. The score of the star model is the sum of responses of the root filter and parts filters, subtracting the displacement cost. Each node in a DPM is labelled with an HOG feature, learned from examples.

By analogy with domain adaptation, we considered the possibility of query expansion for DPM to obtain adapted DPM (ADPM). We first train a standard DPM model for each object category in the training set (i.e., source domain)  $\mathcal{S}$ . We then apply the models on the test set (i.e., target domain)  $\mathcal{T}$ . A confidence set  $\mathcal{C} \subset \mathcal{T}$  is constructed from the test set for training expansion by picking images that match a particular VCM especially well:

$$\mathcal{C} = \{x \in \mathcal{T} | s_1(x) > \theta_1 \wedge s_1(x) - s_2(x) > \theta_2\} \quad (1)$$

with  $s_1(x)$  the highest DPM score, and  $s_2(x)$  the second highest score, and  $\theta_1 \geq \theta_2$  are user-specified parameters to threshold the best score and margin respectively. We found  $\theta_1 = -0.8$  and  $\theta_2 = 0.1$  to be a good trade-off between minimising false positives (5%) and including appropriate number of expanded data (around 580 images in  $\mathcal{C}$ ).

The fully connected multi-labelled graph (MG) model [36] is designed for the cross-depiction problem. It attempts to separate appearance features (contingent on the details of a particular depiction) from the information that characterises an object class without reference to any depiction. Unlike DPM, it comprises a fully connected weighted graph, and has multiple labels per node. Each graph has eight nodes. Weights on nodes can be interpreted as denoting the importance of a node to object class characterisation in a way that is independent of depiction. Weights on arcs are high if the distance between the connected pairs of parts varies little. These weights are learnt using a structural support vector machine [56]. In addition to the weights, each node carries 2 features labels. These



**Fig. 3** Top left: simple shapes learnt from segmentation without supervision. Top right: a hierarchy of shapes derived from an input image. Bottom: a mean graph learnt at head level in the hierarchy, with simple shapes labelling nodes. Edges also connect between layers.

are designed to characterise the appearance of parts in both photographs and artwork (see Section 4 for a justification).

Table 3 compares the classification performance of DPM, ADPM, and MG with the non-structure baseline FV. We can clearly see the benefit when considering the spacial information. Even so, the performance of standard DPM in “train on photo, test on art” pattern significantly drops. However, this performance gap is shortened when the DPM model is re-learned on the expanded set, i.e., ADPM. It demonstrates that the expanded set does capture new information in the target domain and helps to refine the models according to the target domain. The MG alone maintains performance over all train/test patterns. The results suggest that structure and spatial layout is an essential information for recognising an object.

### 3.4 Deep learning

Convolutional neural networks (CNN) [57] have yielded a significant performance boost on image classification. For classification, we follow Crowley and Zisserman [38], encoding images with CNN features, which are then used as input to learn a one-vs-all linear SVM classifier. The CNN parameters are pre-trained from the large ILSVR2013 dataset. We have included results from CNN in Table 2 because they compare so well with the space/structure aware methods. The pre-trained CNN achieved high

**Table 3** Classification using space and structure. Each row is a train (30 image)/test (rest) pattern: Art, Photo, Mixed. Each column is an algorithm. Fisher Vectors [49], the best feature-only classifier, is repeated from Table 1. DPM [55] used a strong spatial layout model, and ADPM is our domain adapted version. Multi-labelled graphs (MG) [36] has a stronger spatial model than DPM, and also has two labels at each node. We have include a deep learning CNN [38] too. Each cell shows the mean of 5 randomised trials. The standard deviation on any column never rises above 2%

Model		FV	DPM	ADPM	MG	CNN
Train	Test	SIFT	HOG	HOG	2×HOG	Learnt
P	P	87	88	—	85	97
M	P	84	85	—	90	96
A	P	66	78	79	83	91
A	A	77	83	—	89	89
M	A	73	80	—	89	87
P	A	47	68	72	83	73

performance when tested on photos. Even so, CNNs exhibit the same fall in performance over the train-on-photo, test-on-art pattern that is seen in the feature based methods.

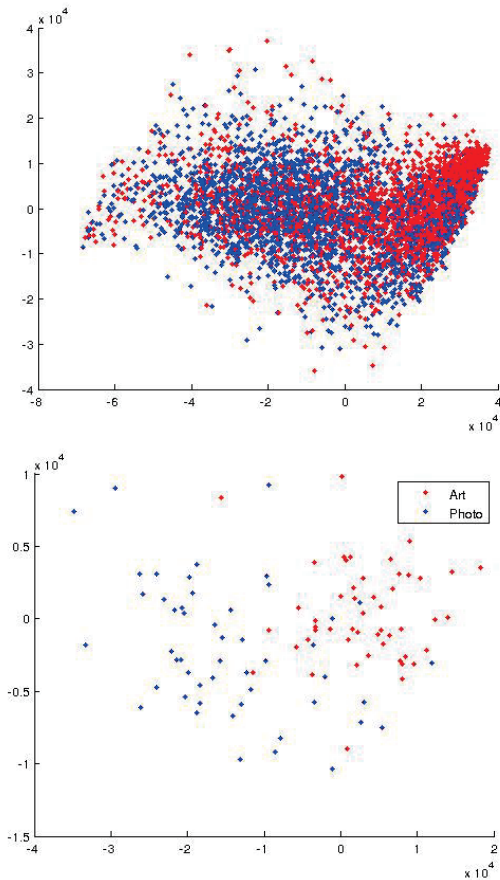
## 4 General discussion

Across all experiments we see the same trend: a fall in performance in any case where art is included. This fall is most marked whenever photographs are used for training and artwork for testing, and is seen in all cases other than the multi-labelled graph (MG) [36].

These observations need an explanation. Intuition suggests that the difference between the low-level images statistics of photographs and artwork is a cause. In particular, it is easy to imagine that the variation in low-level statistics across the gamut of all images is much wider than it is for any one depiction alone (photographs). This intuition is not ours alone, but is shared by others [38], and it remains untested.

A strong hypothesis is possible. Let  $\mathbb{X}$  and  $\mathbb{Y}$  be object classes. Let  $x_P \in \mathbb{X}$  be a photographic instance and  $x_A$  is artwork instance of class  $\mathbb{X}$ . Similarly  $y_P, y_A \in \mathbb{Y}$  are photograph and artwork of class  $\mathbb{Y}$ , respectively. Denote the set of all  $x_P$  by  $X_P$ , meaning the “photo visual object class  $\mathbb{X}$ ”, and likewise for  $X_A, Y_P$ , and  $Y_A$ . Suppose too there is a measure  $d(.,.)$  between each pair of elements in any set. The strong hypothesis is this: the intra-class distance (same domain, different class) is expected to be less than the inter-class distance for (different domain, same class). That is  $d(x_P, x_A) > d(x_P, y_P)$ , photographs are drawings of the same object are more different from each other than photographs of two different objects. Likewise,  $d(x_P, x_A) > d(x_A, y_A)$ , etc. To test this we used raw images Photo-Art-50 as raw input, each scaled to a square image of pixel width 256. We then mapped all the data into a 4 dimensional space using PCA over all the data (which captured most of the eigenenergy). We assumed a K-NN classifier, so that  $X_P$  is represented by the mean, likewise  $X_A$ . The measure,  $d(.,.)$ , is Euclidean distance. We found a fraction 0.67 of all statements of the form  $d(x_P, x_A) > d(x_P, y_P)$ , etc. to be true, which supports the stronger hypothesis. Figure 4





**Fig. 4** Above: each image in Photo-Art-50 plotted in an eigenspace spanning raw images, art in red, photos in blue. Below: The centre of each class in Photo-Art-50, red (art), blue (photo). The images and the cluster centres tend to form two groups: art/photo.

illustrates that for all classes the different domains art/photo tend to separate. This result explain our results above: a density fitted to photographic features alone is over-fitted because it fails to generalise to art-like features, and vice versa. Wu et al. [36] describe feature distributions using more than one centre, and are the most consistent of all descriptions over all recognition tasks on the Photo-Art-50 dataset.

This wide variance in low-level statistics also helps explain the value of spatial information regarding object class identity. So far every method we have experimented that uses some kind of spatial information shows less fall away in the cross-depiction problem; this is true also for Ref. [32]. In this paper we see DPM outperforms BoW, and the MG outperforms DPM. This result is in line with (e.g.) Leordeanu et al. [17] who use the distance between low-level parts (edgelets) as a feature to

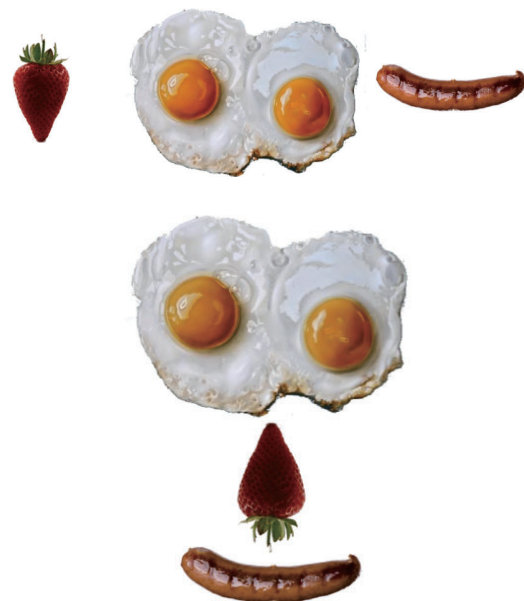
characterise objects and achieve excellent detection results on the PASCAL dataset [41] of photographs; it may be effective too on Photo-Art-50, but this is to be proven.

This empirical data has anecdotal evidence too. The children's drawings in Fig. 1 are clearly people, but have little in common with photographs of people, and not much in common with one another. Consider too Fig. 5 in which the same parts form a face, or not, depending only on the spatial arrangement of the parts. Indeed, artwork from prehistory to the present day, whether produced by a professional or a child, no matter where in the world: the greater majority of it relies on spatial organisation for recognition. It is as if spatial organisation provides a major class, which is refined using features such as shape; but we have no direct evidence for this conjecture.

## 5 Cross-depiction synthesis

Photorealistic image generation is common in computer graphics. Here we focus only on non-photorealistic rendering (NPR) from photographs.

Structure, spatial layout, and shape are all important characteristics in identifying objects regardless of depiction. Equally, they can be used to generate artwork directly from photographs. Consider Fig. 6; it shows a photograph



**Fig. 5** The presence of a face depends on spatial arrangement of parts: above, no face; below, smiling face.

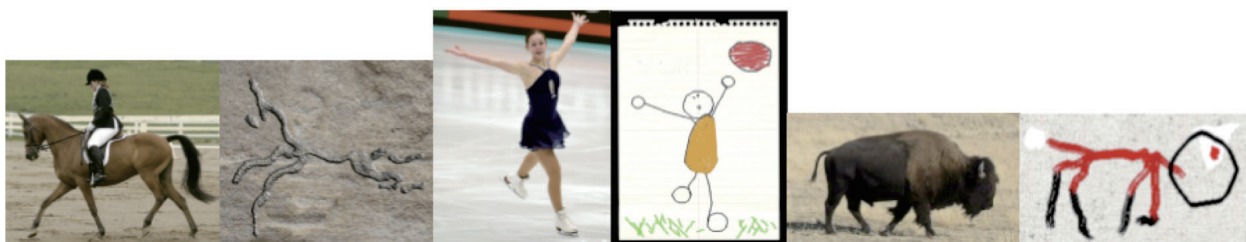


**Fig. 6** Shape abstraction for automated art.

of a bird feeding its young. The photograph has been segmented, and the segments are classified into one of a few qualitative shapes (square, circle, triangle, ...). In the most extreme case just one class (circle) is used. See Ref. [58] for details of the computer graphics algorithm.

It is true that as the degree of abstraction grows the original interpretation of the image becomes harder to maintain; but given too the degree of abstraction in children's drawings, the conclusion is that both the quality and quantity of abstraction are important for recognition. In this case the aim was only to produce a "pretty" image that bears some resemblance to the original. However, simple qualitative shapes of the kind used here can be learned directly from segmentation, as are sufficient to classify scene type (indoor, outdoor, city, ...) at close to state-of-the-art rates [35].

Shape is not the only form of abstraction useful to the production of art, but structure can be used too. Figure 7 shows examples of computer generated art based on rendering structure. The analysis used to obtain the structure is identical to that used by Ref. [37] to classify objects based on weighted graphs alone. In this case the arcs of a graph have been visualised in a non-photorealistic manner, and the shape of parts at nodes have been classified into a qualitative shape; see Ref. [59] for details, which specified the shapes learnt from segmentation by Ref. [35].



**Fig. 7** Structure and shape combine to make art in the style of (left to right) petroglyphs, child art, and Joan Miro.

## 6 Conclusions

It is clear that the same sorts of representations that support abstract image synthesis also support image classification. It seems that synthesis and classification are indeed related, as intuition would have us believe.

The cross-depiction problem pushes at the foundations of computer vision, because it brings in sharp focus the question of how to describe object classes. Given the fact that the same kinds of representations are used both for abstract rendering and for recognition, the conclusion that there is a strong relation between the two is hard to escape. The relation between the cross-depiction problem and image generation is given (strong) anecdotal support by the observation that people draw a mix of what they know and what they see. We can see this in the art of children, and by the fact that when draughting was considered important, by art schools, the tutors had to train students to draw what they see rather than what they know — that is one of the main purposes of life-drawing classes.

Our experimental results show that recognition algorithms premised directly on appearance suffer a fall in performance within the cross-depiction problem, probably because they tacitly assume limited variance of low-level statistics. Rather, they suggest that structure, spatial layout, and shape are all important characteristics in identifying objects regardless of depiction.

For example, DPM outperforms BoW-HOG, even though both use the same low-level features; the MG — with a stronger spatial model — outperforms DPM. This is because, possibly, structure and spatial layout capture the essential form of an object class, with specific appearance relegated to the level of detail. In other words, structure and space are more salient to robust identification that

appearance. Indeed all algorithms we have tested show a significant fall compared to their own peak in performance, when trained on photographs and tested on art; this includes the deep learning methods we have used. The single exception is Ref. [36], which explicitly models a strong structure, and explains appearance details using multiple labels on each node (multiple labels to account for both art and photographic appearance).

The relative importance and the interaction between the descriptors we have identified as important remain an open problem, and does the possibility of other descriptive terms has not been eliminated. A zebra and a horse look largely identical, except for texture.

Deep learning performs very well on classification over Photo-Art-50, but it does exhibit a fall in performance when trained on photographs and tested on art—only the multi-labelled graph [36] and the (lesser performing) graph-with-shapes [34] do not. Also, we have found that when presented with the problem of people detection in a much larger database CNN methods do not rise above a detection rate of 40%. These results make it difficult to conclude that deep learning is a solution to the cross-depiction problem; quite possibly it too suffers from lack of spatial awareness.

In summary, the cross-depiction problem pushes the envelope of computer vision research. It offers significant challenges, which if solved will support new applications in computer graphics and other areas. Modelling visual classes using structure and spatial relations seems to offer a useful way forward; the role of deep learning in the problem is yet to be fully proven in comparison to its own performance in other tasks and when compared to human ability in this difficult challenge.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution License which permits any use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

## References

- [1] Csurka, G.; Dance, C. R.; Fan, L.; Willamowski, J.; Bray, C. Visual categorization with bags of keypoints.

- In: Workshop on Statistical Learning in Computer Vision, ECCV, 1–22, 2004.
- [2] Lazebnik, S.; Schmid, C.; Ponce, J. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, 2169–2178, 2006.
- [3] Russakovsky, O.; Lin, Y.; Yu, K.; Li, F.-F. Object-centric spatial pooling for image classification. *Lecture Notes in Computer Science* 1–15, 2012.
- [4] Gong, B.; Shi, Y.; Sha, F.; Grauman, K. Geodesic flow kernel for unsupervised domain adaptation. In: IEEE Conference on Computer Vision and Pattern Recognition, 2066–2073, 2012.
- [5] Vedaldi, A.; Fulkerson, B. Vlfeat: An open and portable library of computer vision algorithms. In: Proceedings of the international conference on Multimedia, 1469–1472, 2010.
- [6] Gu, C.; Lim, J. J.; Arbelaez, P.; Malik, J. Recognition using regions. In: IEEE Conference on Computer Vision and Pattern Recognition, 1030–1037, 2009.
- [7] Jia, W.; McKenna, S. J. Classifying textile designs using bags of shapes. In: The 20th International Conference on Pattern Recognition, 294–297, 2010.
- [8] Cootes, T. F.; Edwards, G. J.; Taylor, C. J. Active appearance models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 23, No. 6, 681–685, 2001.
- [9] Coughlan, J.; Yuille, A.; English, C.; Snow, D. Efficient deformable template detection and localization without user initialization. *Computer Vision and Image Understanding* Vol. 78, No. 3, 303–319, 2000.
- [10] Crandall, D.; Felzenszwalb, P.; Huttenlocher, D. Spatial priors for part-based recognition using statistical models. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 1, 10–17, 2005.
- [11] Amit, Y.; Trounev, A. Pop: Patchwork of parts models for object recognition. *International Journal of Computer Vision* Vol. 75, No. 2, 267–282, 2007.
- [12] Felzenszwalb, P. F.; Girshick, R. B.; McAllester, D.; Ramanan, D. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 32, No. 9, 1627–1645, 2010.
- [13] Felzenszwalb, P. F.; Huttenlocher, D. P. Pictorial structures for object recognition. *International Journal of Computer Vision* Vol. 61, No. 1, 55–79, 2005.
- [14] Fergus, R.; Perona, P.; Zisserman, A. Object class recognition by unsupervised scale-invariant learning. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, II-264–II-271, 2003.
- [15] Fischler, M. A.; Elschlager, R. A. The representation and matching of pictorial structures. *IEEE Transactions on Computers* Vol. C-22, No. 1, 67–92, 1973.

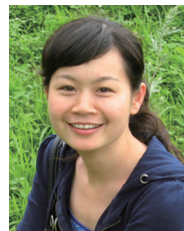
- [16] Leibe, B.; Leonardis, A.; Schiele, B. Robust object detection with interleaved categorization and segmentation. *International Journal of Computer Vision* Vol. 77, Nos. 1–3, 259–289, 2008.
- [17] Leordeanu, M.; Herbert, M.; Sukthankar, R. Beyond local appearance: Category recognition from pairwise interactions of simple features. In: IEEE Conference on Computer Vision and Pattern Recognition, 1–8, 2007.
- [18] Elidan, G.; Heitz, G.; Koller, D. Learning object shape: From drawings to images. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, 2064–2071, 2006.
- [19] Ferrari, V.; Fevrier, L.; Jurie, F.; Schmid, C. Groups of adjacent contour segments for object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 30, No. 1, 36–51, 2008.
- [20] Rom, H.; Medioni, G. Hierarchical decomposition and axial shape description. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 15, No. 10, 973–981, 1993.
- [21] Sundar, H.; Silver, D.; Gagvani, N.; Dickinson, S. Skeleton based shape matching and retrieval. In: Proceedings of the Shape Modeling International, 130–139, 2003.
- [22] Siddiqi, K.; Shokoufandeh, A.; Dickinson, S. J.; Zucker, S. W. Shock graphs and shape matching. *International Journal of Computer Vision* Vol. 35, No. 1, 13–32, 1999.
- [23] Pan, S. J.; Tsang, I. W.; Kwok, J. T.; Yang, Q. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks* Vol. 22, No. 2, 199–210, 2011.
- [24] Gopalan, R.; Li, R.; Chellappa, R. Domain adaptation for object recognition: An unsupervised approach. In: IEEE International Conference on Computer Vision, 999–1006, 2011.
- [25] Fernando, B.; Habrard, A.; Sebban, M.; Tuytelaars, T. Unsupervised visual domain adaptation using subspace alignment. In: IEEE International Conference on Computer Vision, 2960–2967, 2013.
- [26] Crowley, E. J.; Zisserman, A. Of gods and goats: Weakly supervised learning of figurative art. In: British Machine Vision Conference, 2013. Available at <http://www.robots.ox.ac.uk/~vgg/publications/2013/Crowley13/crowley13.pdf>.
- [27] Hu, R.; Collomosse, J. A performance evaluation of gradient field HOG descriptor for sketch based image retrieval. *Computer Vision and Image Understanding* Vol. 117, No. 7, 790–806, 2013.
- [28] Li, Y.; Song, Y.-Z.; Gong, S. Sketch recognition by ensemble matching of structured features. In: Proceedings of the British Machine Vision Conference, 35.1–35.11, 2013.
- [29] Collomosse, J. P.; McNeill, G.; Qian, Y. Storyboard sketches for content based video retrieval. In: IEEE 12th International Conference on Computer Vision, 245–252, 2009.
- [30] Hu, R.; James, S.; Wang, T.; Collomosse, J. Markov random fields for sketch based video retrieval. In: Proceedings of the 3rd ACM conference on International conference on multimedia retrieval, 279–286, 2013.
- [31] Shechtman, E.; Irani, M. Matching local self-similarities across images and videos. In: IEEE Conference on Computer Vision and Pattern Recognition, 1–8, 2007.
- [32] Crowley, E. J.; Zisserman, A. The state of the art: Object retrieval in paintings using discriminative regions. In: British Machine Vision Conference, 2014. Available at <https://www.robots.ox.ac.uk/~vgg/publications/2014/Crowley14/crowley14.pdf>.
- [33] Shrivastava, A.; Malisiewicz, T.; Gupta, A.; Efros, A. A. Data-driven visual similarity for cross-domain image matching. *ACM Transaction of Graphics* Vol. 30, No. 6, Article No. 154, 2011.
- [34] Wu, Q.; Hall, P. Modelling visual objects invariant to depictive style. In: Proceedings of the British Machine Vision Conference, 23.1–23.12, 2013.
- [35] Wu, Q.; Hall, P. Prime shapes in natural images. In: BMCV, 1–12, 2012.
- [36] Wu, Q.; Cai, H.; Hall, P. Learning graphs to model visual objects across different depictive styles. *Lecture Notes in Computer Science* Vol. 8695, 313–328, 2014.
- [37] Xiao, B.; Song Y.-Z.; Hall, P. Learning invariant structure for object identification by using graph methods. *Computer Vision and Image Understanding* Vol. 115, No. 7, 1023–1031, 2011.
- [38] Crowley, E. J.; Zisserman, A. The state of the art: Object retrieval in paintings using discriminative regions. In: British Machine Vision Conference, 2014. Available at <https://www.robots.ox.ac.uk/~vgg/publications/2014/Crowley14/crowley14.pdf>.
- [39] Ginosar, S.; Haas, D.; Brown, T.; Malik, J. Detecting people in cubist art. *Lecture Notes in Computer Science* Vol. 8925, 101–116, 2015.
- [40] BBC. Your paintings dataset. Available at <http://www.bbc.co.uk/arts/yourpaintings/>.
- [41] Everingham, M.; Gool, L. V.; Williams, C. K. I.; Winn, J.; Zisserman, A. The PASCAL visual object classes (voc) challenge. *International Journal of Computer Vision* Vol. 88, No. 2, 303–338, 2010.
- [42] Kyprianidis, J. E.; Collomosse, J.; Wang, T.; Isenberg, T. State of the “art”: A taxonomy of artistic stylization techniques for images and video. *IEEE Transactions on Visualization and Computer Graphics* Vol. 19, No. 5, 866–885, 2013.
- [43] Lowe, D. G. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision* Vol. 60, No. 2, 91–110, 2004.
- [44] Berg, A. C.; Malik, J. Geometric blur for template matching. In: Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 1, I-607–I-614, 2001.

- [45] Chatfield, K.; Philbin, J.; Zisserman, A. Efficient retrieval of deformable shape classes using local self-similarities. In: IEEE 12th International Conference on Computer Vision Workshops, 264–271, 2009.
- [46] Dalal, N.; Triggs, B. Histograms of oriented gradients for human detection. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 1, 886–893, 2005.
- [47] Vedaldi, A.; Zisserman, A. Efficient additive kernels via explicit feature maps. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 34, No. 3, 480–492, 2012.
- [48] Ferrari, V.; Jurie, F.; Schmid, C. From images to shape models for object detection. *International Journal of Computer Vision* Vol. 87, No. 3, 284–303, 2010.
- [49] Perronnin, F.; Sánchez, J.; Mensink, T. Improving the fisher kernel for large-scale image classification. *Lecture Notes in Computer Science* Vol. 6314, 143–156, 2010.
- [50] Hu, R.; Barnard, M.; Collomosse, J. P. Gradient field descriptor for sketch based retrieval and localization. In: The 17th IEEE International Conference on Image Processing, 1025–1028, 2010.
- [51] Gong, B.; Grauman, K.; Sha, F. Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation. In: Proceedings of the International Conference on Machine Learning, 222–230, 2013.
- [52] Saenko, K.; Kulis, B.; Fritz, M.; Darrell, T. Adapting visual category models to new domains. *Lecture Notes in Computer Science* Vol. 6314, 213–226, 2010.
- [53] Song, Y.-Z.; Arbelaez, P.; Hall, P.; Li, C.; Balikai, A. Finding semantic structures in image hierarchies using Laplacian graph energy. *Lecture Notes in Computer Science* Vol. 6314, 694–707, 2010.
- [54] Wu, Q.; Hall, P. Prime shapes in natural images. In: Proceedings of the British Machine Vision Conference, 45.1–45.12, 2012.
- [55] Felzenszwalb, P. F.; Girshick, R. B.; McAllester, D.; Ramanan, D. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 32, No. 9, 1627–1645, 2010.
- [56] Cho, M.; Alahari, K.; Ponce, J. Learning graphs to match. In: Proceedings of the IEEE International Conference on Computer Vision, 25–32, 2013.
- [57] Krizhevsky, A.; Sutskever, I.; Hinton, G. E. ImageNet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25, 1097–1105, 2012.
- [58] Song, Y.-Z.; Pickup, D.; Li, C.; Rosin, P.; Hall, P. Abstract art by shape classification. *IEEE Transactions on Visualization and Computer Graphics* Vol. 19, No. 8, 1252–1263, 2013.

- [59] Hall, P.; Song, Y.-Z. Simple art as abstractions of photographs. In: Proceedings of the Symposium on Computational Aesthetics, 77–85, 2013.



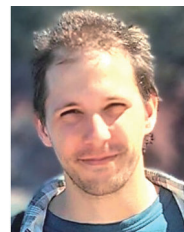
**Peter Hall** is a professor of visual computing at University of Bath and also a director of the Media Technology Research Centre at University of Bath. His research interests cover both computer vision and computer graphics: the relationship between photographs and artwork of all kinds, and 3D model acquisition of complex phenomena from video is of particular interest.



**Hongping Cai** is a postdoctoral researcher in Media Technology Research Centre, University of Bath, UK. She received her Ph.D. degree from National University of Defense Technology, China, in 2010. During her Ph.D., she spent about two years as a visiting student in Centre for Vision, Speech and Signal Processing, University of Surrey, UK. She joint the Centre for Machine Perception, Czech Technical University, Prague, as a post doctor in 2012. Her research interests include cross-depiction classification and detection, texture-less object detection, visual codebook learning, discriminant descriptor learning, and so on.



**Qi Wu** is currently a postdoctoral researcher in Australia Centre for Visual Technologies, University of Adelaide. He received his Ph.D. degree from University of Bath, UK, in 2015. His research interests include cross-depiction object detection and classification, attributes learning, neural networks, image captioning, and so on.



**Tadeo Corradi** is a Ph.D. student at the Mechanical Engineering Department of University of Bath, UK. His research interests include visual and tactile robotics, multi-modal object classification, machine learning, and visuo-tactile integration.

Other papers from this open access journal are available free of charge from <http://www.springer.com/journal/41095>. To submit a manuscript, please go to <https://www.editorialmanager.com/cvmj>.