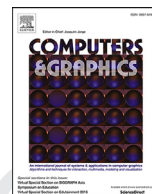




ELSEVIER

Contents lists available at ScienceDirect

Computers &amp; Graphics

journal homepage: [www.elsevier.com/locate/cag](http://www.elsevier.com/locate/cag)

Special Issue on CAD/Graphics 2017

## Joint analysis of shapes and images via deep domain adaptation

 Q1 Zizhao Wu<sup>a,\*</sup>, Yunhui Zhang<sup>a</sup>, Ming Zeng<sup>b</sup>, Feiwei Qin<sup>c</sup>, Yigang Wang<sup>a</sup>
<sup>a</sup>School of Media and Design, Hangzhou Dianzi University, China<sup>b</sup>Software School, Xiamen University, China<sup>c</sup>School of Computer Science, Hangzhou Dianzi University, China

## ARTICLE INFO

## Article history:

Received 15 June 2017

Revised 7 July 2017

Accepted 9 July 2017

Available online xxx

## Keywords:

3D shape recognition

Joint analysis

Cross-modal retrieval

Convolutional neural network

Domain adaptation

## ABSTRACT

3D shapes and 2D images usually contain complementary information for each other, and thus joint analysis of both of them will benefit some problems existing in different domains. Leveraging the connection between 2D images and 3D shapes, it is potential to mine lacking information of one modal from the other. Stemming from this insight, we design and implement a CNN architecture to jointly analyze shapes and images even with few training data guidance. The core of our architecture is a domain adaptation algorithm, which builds up the connection between underlying feature spaces of images and shapes, then aligns and correlates the intrinsic structures therein. The proposed method facilitates the recognition and retrieval tasks. Experiments on the shape recognition tasks show that our approach has superior performance under the difficult setting: zero-shot learning and few-shot learning. We also evaluate our method on the retrieval tasks, and demonstrate the effectiveness of the proposed method.

© 2017 Elsevier Ltd. All rights reserved.

## Q2 1. Introduction

2 With the rapid development of Internet, massive data in multi-  
 3 ple modalities such as images, videos, and 3D models are emerg-  
 4 ing. These heterogeneous data are usually associated to represent  
 5 the same entity. For example, one can generate images, videos, and  
 6 3D models respectively to depict an object in terms of the shape,  
 7 the color, the texture or the motion. The explosive increase of mul-  
 8 timedia data has brought the challenge of how to effectively recog-  
 9 nize, retrieve and organize these resources. Significant efforts have  
 10 been devoted for these tasks, however, most of such efforts handle  
 11 these modalities of data separately, and do not take full advantage  
 12 of the complementary information that exists in different domains.  
 13 This is especially the case in the computer graphics and the com-  
 14 puter vision communities, where the differences in properties of  
 15 viewpoint, lighting, background, occlusion, as well as data repre-  
 16 sentation are most prevalent, hence hinder cross-domain analysis.

17 In recent years, a few works [1–3] have been proposed to ad-  
 18 dress the problem of joint analysis between 3D shapes and 2D  
 19 images. These works have demonstrated the great potential for  
 20 solving some difficult tasks in one domain by exploiting full ad-  
 21 vantage of the complementary information in the other, which in-  
 22 clude cross-view image retrieval, cross-modal retrieval, text based  
 23 shape retrieval, 3D repository and 2D repository filtering, 3D model

alignment, 3D shape recognition, etc. We note that most of these  
 24 works employed the Convolutional Neural Network (CNN) to learn  
 25 the feature vectors for shapes and images, and achieved remark-  
 26 able performance in many fields. These algorithms usually repre-  
 27 sent each 3D shape as a set of rendered images from different  
 28 views around the model, in order to overcome the gap of data rep-  
 29 resentation between 3D shapes and CNN models. However, with  
 30 the exception of Li et al. [3], they all treat the rendered images  
 31 from shapes and 2D images without distinction, which will lead  
 32 the problem of domain bias [4], due to the great discrepancy in  
 33 appearance between them.

34 In this work, we present a novel architecture for joint anal-  
 35 ysis of shapes and images, our method overcomes the discrep-  
 36 ancancy between images of different domains by introducing a do-  
 37 main adaptation algorithm, which leverages knowledge across do-  
 38 mains. Specifically, a joint source and target convolutional neural  
 39 network architecture is introduced to learn the feature represen-  
 40 tations of different domains. Domain adaptation is carried out in  
 41 the learning process, aiming at fusing the feature representations  
 42 of different domains into a shared latent space. The shared fea-  
 43 ture space is semantically meaningful, i.e. data of nearby points  
 44 hold similar semantic information, regardless the modalities they  
 45 belong to.

46 An appealing feature of the proposed algorithm is its capacity of  
 47 leveraging information from one domain to the other, which facil-  
 48 itates the task of dealing with insufficient training data in one do-  
 49 main. As a result, our architecture can be used to few-shot learning  
 50 [5], when a small amount of target labeled data is available from  
 51

\* Corresponding author.

E-mail address: [wuzizhao@163.com](mailto:wuzizhao@163.com) (Z. Wu).

each category, and zero-shot learning [6], when a small amount of target labeled data is available from a subset of the categories. We first evaluated our approach for the recognition on some popular data sets, the results show that our method matches state-of-the-art performance while requiring less training data of the target. Furthermore, we demonstrate that our method has the ability to correctly predict object category labels for unseen categories, i.e. zero-shot classification, by leveraging the knowledge across domains, which is hard for the state-of-the-art recognition methods. We also evaluated the performance of our approach on cross-modal retrieval of shapes and images, and demonstrated its effectiveness.

The main contributions of this paper include (1) we introduce a novel algorithm for joint analysis of 3D shapes and 2D images; and (2) to the best of our knowledge, we are the first to address the few-shot learning and the zero-shot learning problems in the 3D shape recognition field.

The remainder of this paper is organized as follows. We review the related work in Section 2, then present the overview of the proposed algorithm in Section 3. The details of our domain adaptation algorithm is described in Section 4. We show some experimental results on benchmark datasets in Section 5, followed by conclusions and future work in Section 6.

## 2. Related work

Our method is related to prior work on shape descriptors, deep learning, domain adaptation and joint image-shape embedding, which we briefly discuss below.

### 2.1. Shape descriptor

Shape descriptor is an informative representation of the shape, aiming at facilitating tasks such as shape matching, recognition, retrieval, and so on. A large variety of shape descriptors has been developed in the computer graphics community. Most of the earlier shape descriptors focus on geometric properties of the shape such as shape contexts [7,8], shape distribution [9], local diameter [10], volume descriptors [11], spherical harmonics [12], conformal factors [13], Heat Kernel Signature (HKS) [14]. Some view-based descriptors also received widely attention, Cyr and Kimia [15] utilize multi-view projections to identify 3D objects and their poses. Chen et al. [16] propose Light Field Descriptor (LFD), which extracts Fourier descriptors from a set of 2D projections of views.

Instead of designing features according to human prior knowledge, discriminative feature learning provides an alternative way to characterize shapes. This especially benefits from the fast development of deep learning techniques [17], where the learned features lead to significant performance boost in classification and recognition tasks [18,19].

It is only very recent that a few works attempt to tackle 3D shape related problems via deep learning methods, such as classification, recognition and retrieval. Wu et al. [20] work with volumetric representation of 3D shapes, obtaining good results of shape classification on Princeton ModelNet [20]. Zhu et al. [21] utilize Auto-Encoder to learn 3D shape feature with multi-view depth images, leading to accurate 3D shape retrieval. One of the limitations of using 3D volumes as input is the loss in detail when shapes are voxelized. Su et al. [22] propose multi-view CNN(MVCNN) for 3D shape recognition where the features of multiple views are integrated with an extra CNN. Generating shape descriptors based on multiple views can be time-consuming and challenging for real-time retrieval. Bai et al. [23] propose real-time shape retrieval, using GPU acceleration and two inverted files (GIFT).

Our work on 3D shape recognition is similar to MVCNN in that we use deep CNN model to learn shape descriptors. It dif-

fers by the fact that our CNN model considers cross-domain input while MVCNN considers only one modality, which strengthens our method to handle with few-shot learning problem and even zero-shot learning problem.

### 2.2. Convolutional Neural Networks

Recently, Convolutional Neural Networks (CNNs) have been shown to be extremely effective for a variety of visual recognition tasks [18,19,24]. Though many CNN architectures have been proposed, such as AlexNet [18], GoogleNet [25], VGG [26], the network structure of AlexNet remains a popular structure, which consists of five convolutional layers with two fully-connected layers followed by a softmax layer to predict the class label. The network is capable of generating useful feature representations by learning low level features in early convolutional layers and accumulating them to high level semantic features in the latter layers.

There are several deep learning frameworks that efficiently implement the above popular networks, such as Berkeley Caffe [27] and Google Tensorflow [28]. We use Caffe in this paper.

### 2.3. Domain adaptation

Domain adaptation establishes knowledge transfer from the labeled source domain to the unlabeled target domain by exploring domain-invariant structures [29]. Recent studies have shown that deep neural networks can learn more transferable features for domain adaptation [30–32], which produce breakthrough results on some domain adaptation datasets. However these methods solve the problem of domain adaptation within the same modality. It is unclear how this can be done when moving across modalities. To address this issue, some notable approaches focus on the problem of jointly embedding or learning representations from multiple modalities into a shared feature space to improve learning [33] or enabling zero-shot learning [34,35].

Our work is primarily motivated by Tzeng et al. [36], that introduces a deep CNN model for the domain adaptation problem. We introduce this architecture to facilitate our task for joint analysis of shapes and images, and further make some optimizations to the original model.

### 2.4. Joint image-shape analysis

Multi-modal feature learning has been researched thoroughly over the past years [33], whereas only few works have addressed the problem in the computer graphics field in recent years. Herzog et al. [1] suggest a new method for structuring multi-modal representations of shapes and keywords, and adds images and sketches to the mix. This method builds the embedding mainly upon the class co-relation and hand-crafted descriptors. Hueting et al. [2] introduce a system for joint images and shapes processing to

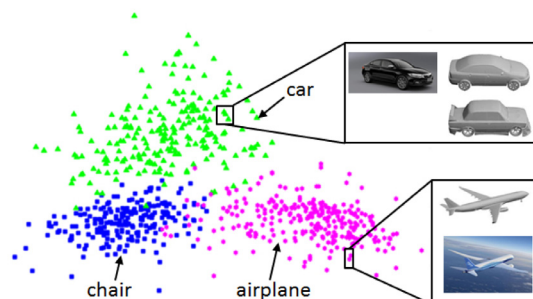


Fig. 1. Our method on joint analysis of shapes and images, and learn to obtain a joint representation in a shared latent space.

Download English Version:

<https://daneshyari.com/en/article/6876861>

Download Persian Version:

<https://daneshyari.com/article/6876861>

[Daneshyari.com](https://daneshyari.com)