

CNN-PS: CNN-Based Photometric
Stereo for General Non-convex Surfaces
Satoshi Ikehata(B)
National Institute of Informatics, Tokyo, Japan
sikehata@nii.ac.jp

Abstract. Most conventional photometric stereo algorithms inversely solve a BRDF-based image formation model. However, the actual imaging process is often far more complex due to the global light transport on the non-convex surfaces. This paper presents a photometric stereo network that directly learns relationships between the photometric stereo input and surface normals of a scene. For handling unordered, arbitrary number of input images, we merge all the input data to the intermediate representation called observation map that has a fixed shape, is able to be fed into a CNN. To improve both training and prediction, we take into account the rotational pseudo-invariance of the observation map that is derived from the isotropic constraint. For training the network, we create a synthetic photometric stereo dataset that is generated by a physics-based renderer, therefore the global light transport is considered. Our experimental results on both synthetic and real datasets show that our method outperforms conventional BRDF-based photometric stereo algorithms especially when scenes are highly non-convex.

Keywords: Photometric stereo
·Convolutional neural networks

1

Dynamic Conditional Networks
for Few-Shot Learning

Fang Zhao¹, Jian Zhao^{1,2}(B), Shuicheng Yan^{1,3}, and Jiashi Feng¹
¹National University of Singapore, Singapore, Singapore
{elezhf,eleyans,elefjia}@nus.edu.sg
²National University of Defense Technology, Hunan, China
³Qihoo 360 AI Institute, Beijing, China
zhaojian90@u.nus.edu
<https://zhaoj9014.github.io/>

Abstract. This paper proposes a novel Dynamic Conditional Convolutional Network (DCCN) to handle conditional few-shot learning, i.e., only a few training samples are available for each condition. DCCN consists of dual subnets: DyConvNet contains a dynamic convolutional layer with a bank of basis filters; CondiNet predicts a set of adaptive weights from conditional inputs to linearly combine the basis filters. In this manner, a specific convolutional kernel can be dynamically obtained for each conditional input. The filter bank is shared between all conditions thus only a low-dimension weight vector needs to be learned. This significantly facilitates the parameter learning across different conditions when training data are limited. We evaluate DCCN on four tasks which can be formulated as conditional model learning, including specific object counting, multi-modal image classification, phrase grounding and identity based face generation. Extensive experiments demonstrate the superiority of the proposed model in the conditional few-shot learning setting.

Keywords: Conditional model
·Few-shot learning ·Deep learning
Dynamic convolution ·Filter bank

1

Deep Factorised Inverse-Sketching

Kaiyue Pang¹(B), Da Li¹, Jifei Song¹, Yi-Zhe Song¹, Tao Xiang¹,
and Timothy M. Hospedales^{1,2}
¹SketchX, Queen Mary University of London, London, UK
{kaiyue.pang,da.li,yizhe.song,j.song,t.xiang}@qmul.ac.uk
²The University of Edinburgh, Edinburgh, UK

t.hospedales@ed.ac.uk

Abstract. Modelling human free-hand sketches has become topical recently, driven by practical applications such as fine-grained sketch based image retrieval (FG-SBIR). Sketches are clearly related to photoedge-maps, but a human free-hand sketch of a photo is not simply a clean rendering of that photo's edge map. Instead there is a fundamental process of abstraction and iconic rendering, where overall geometry is warped and salient details are selectively included. In this paper we study this sketching process and attempt to invert it. We model this inversion by translating iconic free-hand sketches to contours that resemble more geometrically realistic projections of object boundaries, and separately factorise out the salient added details. This factorised re-representation makes it easier to match a free-hand sketch to a photo instance of an object. Specifically, we propose a novel unsupervised image style transfer model based on enforcing a cyclic embedding consistency constraint. A deep FG-SBIR model is then formulated to accommodate complementary discriminative detail from each factorised sketch for better matching with the corresponding photo. Our method is evaluated both qualitatively and quantitatively to demonstrate its superiority over a number of state-of-the-art alternatives for style transfer and FG-SBIR.

1

Look Deeper into Depth: Monocular
Depth Estimation with Semantic Booster
and Attention-Driven Loss

Jianbo Jiao^{1,2(B)}, Ying Cao¹, Yibing Song³, and Rynson Lau¹

¹City University of Hong Kong, Kowloon, Hong Kong SAR

jianbjiao2-c@my.cityu.edu.hk, caoying59@gmail.com,

dynamicstevenson@gmail.com

²University of Illinois at Urbana-Champaign, Urbana, USA

³Tencent AI Lab, Shenzhen, China

rynson.lau@cityu.edu.hk

Abstract. Monocular depth estimation benefits greatly from learning based techniques. By studying the training data, we observe that the per-pixel depth values in existing datasets typically exhibit a long-tailed distribution. However, most previous approaches treat all the regions in the training data equally regardless of the imbalanced depth distribution, which restricts the model performance particularly on distant depth regions. In this paper, we investigate the long tail property and delve deeper into the distant depth regions (i.e. the tail part) to propose an attention-driven loss for the network supervision. In addition, to better leverage the semantic information for monocular depth estimation, we propose a synergy network to automatically learn the information sharing strategies between the two tasks. With the proposed attention-driven loss and synergy network, the depth estimation and semantic labeling tasks can be mutually improved. Experiments on the challenging indoor datasets show that the proposed approach achieves state-of-the-art performance on both monocular depth estimation and semantic labeling tasks.

Keywords: Monocular depth

· Semantic labeling · Attention loss

1

Summarizing First-Person Videos from
Third Persons' Points of Views

Hsuan-I Hol¹, Wei-Chen Chiu², and Yu-Chiang Frank Wang^{1(B)}

¹Department of Electrical Engineering, National Taiwan University, Taiwan, Republic of China

{b01901029, ycwang}@ntu.edu.tw

²Department of Computer Science, National Chiao Tung University, Taiwan, Republic of China

walon@cs.nctu.edu.tw

Abstract. Video highlight or summarization is among interesting topics in computer vision, which benefits a variety of applications like viewing, searching, or storage. However, most existing studies rely on training data of third-person videos, which cannot easily generalize to highlight the first-person ones. With the goal of deriving an effective model to summarize first-person videos, we propose a novel deep neural network architecture for describing and discriminating vital spatiotemporal information across videos with different points of view. Our proposed model is realized in a semi-supervised setting, in which fully annotated third-person videos, unlabeled first-person videos, and a small number of annotated first-person ones are presented during training. In our experiments, qualitative and quantitative evaluations on both benchmarks and our collected first-person video datasets are presented.

Keywords: Video summarization

•First-person vision

Transfer learning •Metric learning

1

Learning Single-View 3D Reconstruction

with Limited Pose Supervision

Guandao Yang^{1(B)}, Yin Cui^{1,2}, Serge Belongie^{1,2}, and Bharath Hariharan¹

¹Department of Computer Science, Cornell University, Ithaca, USA

gy46@cornell.edu

²Cornell Tech, New York, USA

Abstract. It is expensive to label images with 3D structure or precise camera pose. Yet, this is precisely the kind of annotation required to train single-view 3D reconstruction models. In contrast, unlabeled images or images with just category labels are easy to acquire, but few current models can use this weak supervision. We present a unified framework that can combine both types of supervision: a small amount of camera pose annotations are used to enforce pose-invariance and view-point consistency, and unlabeled images combined with an adversarial loss are used to enforce the realism of rendered, generated models. We use this unified framework to measure the impact of each form of supervision in three paradigms: semi-supervised, multi-task, and transfer learning. We show that with a combination of these ideas, we can train single-view reconstruction models that improve up to 7 points in performance (AP) when using only 1% pose annotated training data.

Keywords: Single-image 3D-reconstruction

•Few-shot learning

GANs

1

Weakly- and Semi-supervised Panoptic

Segmentation

Qizhu Li, Anurag Arnab^(B), and Philip H.S. Torr

University of Oxford, Oxford, UK

{liqizhu,aarnab,phst}@robots.ox.ac.uk

Abstract. We present a weakly supervised model that jointly performs both semantic- and instance-segmentation – a particularly relevant problem given the substantial cost of obtaining pixel-perfect annotation for these tasks. In contrast to many popular instance segmentation approaches based on object detectors, our method does not predict any overlapping instances. Moreover, we are able to segment both “thing” and “stuff” classes, and thus explain all the pixels in the image. “Thing” classes are weakly-supervised with bounding boxes, and “stuff” with image-level tags. We obtain state-of-the-art results on Pascal VOC, for both full and weak supervision (which achieves about 95% of fully-supervised performance). Furthermore, we present the first weakly-sup

ervised results on Cityscapes for both semantic- and instance-segmentation. Finally, we use our weakly supervised framework to analyse the relationship between annotation quality and predictive performance, which is of interest to dataset creators.

Keywords: Weak supervision

•Instance segmentation

Semantic segmentation •Scene understanding

1

Making Deep Heatmaps Robust to Partial

Occlusions for 3D Object Pose Estimation

Markus Oberweger^{1(B)}, Mahdi Rad¹, and Vincent Lepetit^{2,1}

¹Institute for Computer Graphics and Vision, Graz University of Technology, Graz, Austria

{oberweger,rad,lepetit }@icg.tugraz.at

²Laboratoire Bordelais de Recherche en Informatique, Université de Bordeaux, Bordeaux, France

Abstract. We introduce a novel method for robust and accurate 3D object pose estimation from a single color image under large occlusions. Following recent approaches, we first predict the 2D projections of 3D points related to the target object and then compute the 3D pose from these correspondences using a geometric method. Unfortunately, as the results of our experiments show, predicting these 2D projections using a regular CNN or a Convolutional Pose Machine is highly sensitive to partial occlusions, even when these methods are trained with partially occluded examples. Our solution is to predict heatmaps from multiple small patches independently and to accumulate the results to obtain accurate and robust predictions. Training subsequently becomes challenging because patches with similar appearances but different positions on the object correspond to different heatmaps. However, we provide a simple yet effective solution to deal with such ambiguities. We show that our approach outperforms existing methods on two challenging datasets: The Occluded LineMOD dataset and the YCB-Video dataset, both exhibiting cluttered scenes with highly occluded objects.

Keywords: 3D object pose estimation

•Heatmaps •Occlusions

1

Deep Co-Training for Semi-Supervised

Image Recognition

Siyuan Qiao^{1(B)}, Wei Shen^{1,2}, Zhishuai Zhang¹, Bowen Geng³, and Alan Yuille¹

¹Johns Hopkins University, Baltimore, USA

siyuan.qiao@jhu.edu

²Shanghai University, Shanghai, China

³Hikvision Research Institute, Hangzhou, China

Abstract. In this paper, we study the problem of semi-supervised image recognition, which is to learn classifiers using both labeled and unlabeled images. We present Deep Co-Training, a deep learning based method inspired by the Co-Training framework. The original Co-Training learns two classifiers on two views which are data from different sources that describe the same instances. To extend this concept to deep learning, Deep Co-Training trains multiple deep neural networks to be the different views and exploits adversarial examples to encourage view difference, in order to prevent the networks from collapsing into each other. As a result, the co-trained networks provide different and complementary information about the data, which is necessary for the Co-Training framework to achieve good results. We test our method on SVHN, CIFAR-10/100 and ImageNet datasets, and our method outperforms the previous state-of-the-art methods by a large margin.

Keywords: Co-Training

•Deep networks •Semi-supervised learning

1

Visual Coreference Resolution in Visual Dialog Using Neural Module Networks

Satwik Kottur^{1,2(B)}, Josée M. F. Moura², Devi Parikh^{1,3}, Dhruv Batra^{1,3},
and Marcus Rohrbach¹

¹Facebook AI Research, Menlo Park, USA

²Carnegie Mellon University, Pittsburgh, USA

skottur@andrew.cmu.edu

³Georgia Institute of Technology, Atlanta, USA

Abstract. Visual dialog entails answering a series of questions grounded in an image, using dialog history as context. In addition to the challenges found in visual question answering (VQA), which can be seen as one-round dialog, visual dialog encompasses several more. We focus on one such problem called visual coreference resolution that involves determining which words, typically noun phrases and pronouns, co-refer to the same entity/object instance in an image. This is crucial, especially for pronouns (e.g., 'it'), as the dialog agent must first link it to a previous coreference (e.g., 'boat'), and only then can rely on the visual grounding of the coreference 'boat' to reason about the pronoun 'it'. Prior work (in visual dialog) models visual coreference resolution either (a) implicitly via a memory network over history, or (b) at a coarse level for the entire question; and not explicitly at a phrase level of granularity. In this work, we propose a neural module network architecture for visual dialog by introducing two novel modules—Refer and Exclude—that perform explicit, grounded, coreference resolution at a finer word level. We demonstrate the effectiveness of our model on MNIST Dialog, a visually simple yet coreference-wise complex dataset, by achieving near perfect accuracy, and on VisDial, a large and challenging visual dialog dataset on real images, where our model outperforms other approaches, and is more interpretable, grounded, and consistent qualitatively.

1

Learning Blind Video Temporal Consistency

Wei-Sheng Lai^{1(B)}, Jia-Bin Huang², Oliver Wang³, Eli Shechtman³,
Ersin Yumer⁴, and Ming-Hsuan Yang^{1,5}

¹UC Merced, Merced, USA

wlai24@ucmerced.edu

²Virginia Tech, Blacksburg, USA

³Adobe Research, Seattle, USA

⁴Argo AI, Mountain View, USA

⁵Google Cloud AI, Sunnyvale, USA

Abstract. Applying image processing algorithms independently to each frame of a video often leads to undesired inconsistent results overtime. Developing temporally consistent video-based extensions, however, requires domain knowledge for individual tasks and is unable to generalize to other applications. In this paper, we present an efficient approach based on a deep recurrent network for enforcing temporal consistency in a video. Our method takes the original and per-frame processed videos as inputs to produce a temporally consistent video. Consequently, our approach is agnostic to specific image processing algorithms applied to the original video. We train the proposed network by minimizing both short-term and long-term temporal losses as well as a perceptual loss to strike a balance between temporal coherence and perceptual similarity with the processed frames. At test time, our model does not require computing optical flow and thus achieves real-time speed even for high-resolution videos. We show that our single model can handle multiple and unseen

tasks, including but not limited to artistic style transfer, enhancement, colorization, image-to-image translation and intrinsic image decomposition. Extensive objective evaluation and subject study demonstrate that the proposed approach performs favorably against the state-of-the-art methods on various types of videos.

1

Salient Objects in Clutter:

Bringing Salient Object Detection
to the Foreground

Deng-Ping Fan¹, Ming-Ming Cheng^{1(B)}, Jiang-Jiang Li¹,
Shang-Hua Gao¹, Qibin Hou¹, and Ali Borji²

¹College of Computer Science, Nankai University, Tianjin, China
cmm@nankai.edu.cn

²CRCV, University of Central Florida, Orlando, Florida, US

<http://mmcheng.net/SOCBenchmark/>

Abstract. We provide a comprehensive evaluation of salient object detection (SOD) models. Our analysis identifies a serious design bias of existing SOD datasets which assumes that each image contains at least one clearly outstanding salient object in low clutter. The design bias has led to a saturated high performance for state-of-the-art SOD models when evaluated on existing datasets. The models, however, still perform far from being satisfactory when applied to real-world daily scenes. Based on our analyses, we first identify 7 crucial aspects that a comprehensive and balanced dataset should fulfill. Then, we propose a new high quality dataset and update the previous saliency benchmark.

Specifically, our SOC (Salient Objects in Clutter) dataset, includes images with salient and non-salient objects from daily object categories. Beyond object category annotations, each salient image is accompanied by attributes that reflect common challenges in real-world scenes. Finally, we report attribute-based performance assessment on our dataset.

Keywords: Salient object detection

· Saliency benchmark

Dataset · Attribute

1

Gray-Box Adversarial Training

B. S. Vivek^(B), Konda Reddy Mopuri, and R. Venkatesh Babu
Indian Institute of Science, Bangalore, India

svivek@iisc.ac.in, kondamopuri@iisc.ac.in, venky@iisc.ac.in

Abstract. Adversarial samples are perturbed inputs crafted to mislead the machine learning systems. A training mechanism, called adversarial training, which presents adversarial samples along with clean samples has been introduced to learn robust models. In order to scale adversarial training for large datasets, these perturbations can only be crafted using fast and simple methods (e.g., gradient ascent). However, it is shown that adversarial training converges to a degenerate minimum, where the model appears to be robust by generating weaker adversaries. As a result, the models are vulnerable to simple black-box attacks.

In this paper we, (i) demonstrate the shortcomings of existing evaluation policy, (ii) introduce novel variants of white-box and black-box attacks, dubbed "gray-box adversarial attacks" based on which we propose novel evaluation method to assess the robustness of the learned models, and (iii) propose a novel variant of adversarial training, named "Gray-box Adversarial Training" that uses intermediate versions of the models to seed the adversaries. Experimental evaluation demonstrates that the models trained using our method exhibit better robustness compared to both undefended and adversarially trained models.

Keywords: Adversarial perturbations

Attacks on machine learning models

•Adversarial training

Robust machine learning models

1

Visual Question Answering as a Meta Learning Task

Damien Teney(B)and Anton van den Hengel

Australian Institute for Machine Learning, University of Adelaide,
Adelaide, Australia

{damien.teney,anton.vandenhengel }@adelaide.edu.au

Abstract. The predominant approach to Visual Question Answering (VQA) demands that the model represents within its weights all of the information required to answer any question about any image. Learning this information from any real training set seems unlikely, and representing it in a reasonable number of weights doubly so. We propose instead to approach VQA as a meta learning task, thus separating the question answering method from the information required. At test time, the method is provided with a support set of example questions/answers, over which it reasons to resolve the given question. The support set is not fixed and can be extended without retraining, thereby expanding the capabilities of the model. To exploit this dynamically provided information, we adapt a state-of-the-art VQA model with two techniques from the recent meta learning literature, namely prototypical networks and meta networks. Experiments demonstrate the capability of the system to learn to produce completely novel answers (i.e. never seen during training) from examples provided at test time. In comparison to the existing state of the art, the proposed method produces qualitatively distinct results with higher recall of rare answers, and a better sample efficiency that allows training with little initial data. More importantly, it represents an important step towards vision-and-language methods that can learn and reason on-the-fly.

The task of Visual Question Answering (VQA) demands that an agent correctly answer a previously unseen question about a previously unseen image. The fact that neither the question nor the image is specified until test time means that the agent must embody most of the achievements of Computer Vision and Natural Language Processing, and many of those of Artificial Intelligence.

VQA is typically framed in a purely supervised learning setting. A large training set of example questions, images, and their correct answers is used to train a method to map a question and image to scores over a predetermined, fixed vocabulary of possible answers using the maximum likelihood [39]. This approach has inherent scalability issues, as it attempts to represent all world Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-01267-0_14) contains supplementary material, which is available to authorized users.

c/circlecopyrtSpringer Nature Switzerland AG 2018

V. Ferrari et al. (Eds.): ECCV 2018, LNCS 11219, pp. 229–245, 2018.https://doi.org/10.1007/978-3-030-01267-0_14

_1

On Online Evaluation of Vision-Based Driving Models

Felipe Codevilla¹, Antonio M. López¹, Vladlen Koltun²,
and Alexey Dosovitskiy³(B)

¹Computer Vision Center, Universitat Autònoma de Barcelona, Barcelona, Spain
fcodevilla@cvc.uab.es

²Intel Labs, Santa Clara, USA

³Intel Labs, Munich, Germany

adosovitskiy@gmail.com

Abstract. Autonomous driving models should ideally be evaluated by deploying them on a fleet of physical vehicles in the real world. Unfortunately,

this approach is not practical for the vast majority of researchers. An attractive alternative is to evaluate models offline, on a pre-collected validation dataset with ground truth annotation. In this paper, we investigate the relation between various online and offline metrics for evaluation of autonomous driving models. We find that offline prediction errors are not necessarily correlated with driving quality, and two models with identical prediction error can differ dramatically in their driving performance. We show that the correlation of offline evaluation with driving quality can be significantly improved by selecting an appropriate validation dataset and suitable offline metrics.

Keywords: Autonomous driving

•Deep learning

1

Visual Psychophysics for Making Face

Recognition Algorithms More Explainable

Brandon Richard Webster^{1(B)}, So Yon Kwon², Christopher Clarizio¹,

Samuel E. Anthony^{2,3}, and Walter J. Scheirer¹

¹University of Notre Dame, Notre Dame, IN 46556, USA

bricharl@nd.edu

²Perceptive Automata, Inc., Somerville, USA

³Harvard University, Cambridge, MA 02138, USA

Abstract. Scientists fields that are interested in faces have developed their own sets of concepts and procedures for understanding how a target model system (be it a person or algorithm) perceives a face under varying conditions. In computer vision, this has largely been in the form of dataset evaluation for recognition tasks where summary statistics are used to measure progress. While aggregate performance has continued to improve, understanding individual causes of failure has been difficult, as it is not always clear why a particular face fails to be recognized, or why an impostor is recognized by an algorithm. Importantly, other fields studying vision have addressed this via the use of visual psychophysics: the controlled manipulation of stimuli and careful study of the responses they evoke in a model system. In this paper, we suggest that visual psychophysics is a viable methodology for making face recognition algorithms more explainable. A comprehensive set of procedures is developed for assessing face recognition algorithm behavior, which is then deployed over state-of-the-art convolutional neural networks and more basic, yet still widely used, shallow and handcrafted feature-based approaches.

Keywords: Face recognition

•Biometrics •Explainable AI

Visual psychophysics •Biometric menagerie

1

Conditional Prior Networks for Optical

Flow

Yanchao Yang^(B) and Stefano Soatto

UCLA Vision Lab, University of California, Los Angeles, CA 90095, USA

{yanchao.yang,soatto}@cs.ucla.edu

Abstract. Classical computation of optical flow involves generic priors (regularizers) that capture rudimentary statistics of images, but not long-range correlations or semantics. On the other hand, fully supervised methods learn the regularity in the annotated data, without explicit regularization and with the risk of overfitting. We seek to learn richer priors on the set of possible flows that are statistically compatible with an image. Once the prior is learned in a supervised fashion, one can easily learn the full map to infer optical flow directly from two or more images, without any need for (additional) supervision. We introduce a novel architecture, called Conditional Prior Network (CPN), and show how to train it to yield a conditional prior. When used in conjunction with a sim

ple optical flow architecture, the CPN beats all variational methods and all unsupervised learning-based ones using the same data term. It performs comparably to fully supervised ones, that however are fine-tuned to a particular dataset. Our method, on the other hand, performs well even when transferred between datasets. Code is available at: <https://github.com/YanchaoYang/Conditional-Prior-Networks> .

1

Robust Optical Flow in Rainy Scenes

Ruoteng Lil(B), Robby T. Tan^{1,2}, and Loong-Fah Cheong¹

¹National University of Singapore, Singapore, Singapore

liruoteng@gmail.com

²Yale-NUS College, Singapore, Singapore

Abstract. Optical flow estimation in rainy scenes is challenging due to degradation caused by rain streaks and rain accumulation, where the latter refers to the poor visibility of remote scenes due to intense rainfall. To resolve the problem, we introduce a residue channel, a single channel (gray) image that is free from rain, and its colored version, a colored-residue image. We propose to utilize these two rain-free images in computing optical flow. To deal with the loss of contrast and the attendant sensitivity to noise, we decompose each of the input images into a piecewise-smooth structure layer and a high-frequency fine-detail texture layer. We combine the colored-residue images and structure layers in a unified objective function, so that the estimation of optical flow can be more robust. Results on both synthetic and real images show that our algorithm outperforms existing methods on different types of rain sequences. To our knowledge, this is the first optical flow method specifically dealing with rain. We also provide an optical flow dataset consisting of both synthetic and real rain images.

Keywords: Optical flow

·Rain·Decomposition ·Residue channel

1

Rethinking Spatiotemporal Feature

Learning: Speed-Accuracy Trade-offs

in Video Classification

Saining Xie^{1,2}, Chen Sun¹(B), Jonathan Huang¹, Zhuowen Tu^{1,2},

and Kevin Murphy¹

¹Google Research, Mountain View, USA

chensun@google.com

²University of California San Diego, San Diego, USA

Abstract. Despite the steady progress in video analysis led by the adoption of convolutional neural networks (CNNs), the relative improvement has been less drastic as that in 2D static image classification. Three main challenges exist including spatial (image) feature representation, temporal information representation, and model/computation complexity. It was recently shown by Carreira and Zisserman that 3D CNNs, inherited from 2D networks and pretrained on ImageNet, could be a promising way for spatial and temporal representation learning. However, as for model/computation complexity, 3D CNNs are much more expensive than 2D CNNs and prone to overfit. We seek a balance between speed and accuracy by building an effective and efficient video classification system through systematic exploration of critical network design choices. In particular, we show that it is possible to replace many of the 3D convolutions by low-cost 2D convolutions. Rather surprisingly, best result (in both speed and accuracy) is achieved when replacing the 3D convolutions at the bottom of the network, suggesting that temporal representation learning on high-level "semantic" features is more useful. Our conclusion generalizes to datasets with very different properties. When combined with several other cost-effective designs including separable spa-

tial/temporal convolution and feature gating, our system results in an effective video classification system that produces very competitive results on several action classification benchmarks (Kinetics, Something-something, UCF101 and HMDB), as well as two action detection (localization) benchmarks (JHMDB and UCF101-24).

1

Variational Wasserstein Clustering

Liang Mil(B), Wen Zhang¹, Xianfeng Gu², and Yalin Wang¹

¹Arizona State University, Tempe, USA

{liangmi,wzhan139}@asu.edu

ylwang@asu.edu

²Stony Brook University, Stony Brook, USA

gu@cs.stonybrook.edu

Abstract. We propose a new clustering method based on optimal transportation. We discuss the connection between optimal transportation and k-means clustering, solve optimal transportation with the variational principle, and investigate the use of power diagrams as transportation plans for aggregating arbitrary domains into a fixed number of clusters.

We drive cluster centroids through the target domain while maintaining the minimum clustering energy by adjusting the power diagram. Thus, we simultaneously pursue clustering and the Wasserstein distance between the centroids and the target domain, resulting in a measure-preserving mapping. We demonstrate the use of our method in domain adaptation, remeshing, and learning representations on synthetic and real data.

Keywords: Clustering

·Discrete distribution ·K-means

Measure preserving ·Optimal transportation ·Wasserstein distance

1

Show, Tell and Discriminate: Image

Captioning by Self-retrieval with Partially

Labeled Data

Xihui Liu¹, Hongsheng Li¹(B), Jing Shao², Dapeng Chen¹,

and Xiaogang Wang¹

¹The Chinese University of Hong Kong, Hong Kong, China

{hsli,xgwan}@ee.cuhk.edu.hk, xihui-liu@link.cuhk.edu.hk,

dpchen@cuhk.edu.hk

²SenseTime Research, Hong Kong, China

shaoming@sensetime.com

Abstract. The aim of image captioning is to generate captions by machine to describe image contents. Despite many efforts, generating discriminative captions for images remains non-trivial. Most traditional approaches imitate the language structure patterns, thus tend to fall into a stereotype of replicating frequent phrases or sentences and neglect unique aspects of each image. In this work, we propose an image captioning framework with a self-retrieval module as training guidance, which encourages generating discriminative captions. It brings unique advantages: (1) the self-retrieval guidance can act as a metric and an evaluator of caption discriminativeness to assure the quality of generated captions. (2) The correspondence between generated captions and images are naturally incorporated in the generation process without human annotations, and hence our approach could utilize a large amount of unlabeled images to boost captioning performance with no additional annotations. We demonstrate the effectiveness of the proposed retrieval-guided method on COCO and Flickr30k captioning datasets, and show its superior captioning performance with more discriminative captions.

Keywords: Image captioning

·Language and vision

Text-image retrieval

1

Contour Knowledge Transfer for Salient Object Detection

Xin Li¹, Fanyang Yang¹, Hong Cheng¹, Wei Liu¹, and Dinggang Shen²

¹University of Electronic Science and Technology of China, Chengdu 611731, China
xinli_uestc@hotmail.com, fanyang_uestc@hotmail.com, hcheng@uestc.edu.cn

²Department of Radiology and BRIC, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA

dgshen@med.unc.edu

Abstract. In recent years, deep Convolutional Neural Networks (CNNs) have broken all records in salient object detection. However, training such a deep model requires a large amount of manual annotations. Our goal is to overcome this limitation by automatically converting an existing deep contour detection model into a salient object detection model without using any manual salient object masks. For this purpose, we have created a deep network architecture, namely Contour-to-Saliency Network (C2S-Net), by grafting a new branch onto a well-trained contour detection network. Therefore, our C2S-Net has two branches for performing two different tasks: (1) predicting contours with the original contour branch, and (2) estimating per-pixel saliency score of each image with the newly-added saliency branch. To bridge the gap between these two tasks, we further propose a contour-to-saliency transferring method to automatically generate salient object masks which can be used to train the saliency branch from outputs of the contour branch. Finally, we introduce a novel alternating training pipeline to gradually update the network parameters. In this scheme, the contour branch generates saliency masks for training the saliency branch, while the saliency branch, in turn, feeds back saliency knowledge in the form of saliency-aware contour labels, for fine-tuning the contour branch. The proposed method achieves state-of-the-art performance on five well-known benchmarks, outperforming existing fully supervised methods while also maintaining high efficiency.

Keywords: Saliency detection

•Deep learning •Transfer learning

1

Learning Category-Specific Mesh

Reconstruction from Image Collections

Angjoo Kanazawa^(B), Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik

University of California, Berkeley, USA

{kanazawa,shubhtuls,efros,malik}@eecs.berkeley.edu

Abstract. We present a learning framework for recovering the 3D shape, camera, and texture of an object from a single image. The shape is represented as a deformable 3D mesh model of an object category where a shape is parameterized by a learned mean shape and per-instance predicted deformation. Our approach allows leveraging an annotated image collection for training, where the deformable model and the 3D prediction mechanism are learned without relying on ground-truth 3D or multi-view supervision. Our representation enables us to go beyond existing 3D prediction approaches by incorporating texture inference as prediction of an image in a canonical appearance space. Additionally, we show that semantic key points can be easily associated with the predicted shapes.

We present qualitative and quantitative results of our approach on CUB and PASCAL3D datasets and show that we can learn to predict diverse shapes and textures across objects using only annotated image collections.

The project website can be found at <https://akanazawa.github.io/cmr/>.

1

Learning to Forecast and Refine Residual

Motion for Image-to-Video Generation

Long Zhao^{1(B)}, Xi Peng², Yu Tian¹, Mubbasir Kapadia¹,
and Dimitris Metaxas¹

¹Rutgers University, Piscataway, USA

{lzh311,yt219,mk1353,dnm }@cs.rutgers.edu

²Binghamton University, Binghamton, USA

xpeng@binghamton.edu

Abstract. We consider the problem of image-to-video translation, where an input image is translated into an output video containing motions of a single object. Recent methods for such problems typically train transformation networks to generate future frames conditioned on the structure sequence. Parallel work has shown that short high-quality motions can be generated by spatiotemporal generative networks that leverage temporal knowledge from the training data. We combine the benefits of both approaches and propose a two-stage generation framework where videos are generated from structures and then refined by temporal signals. To model motions more efficiently, we train networks to learn residual motion between the current and future frames, which avoids learning motion-irrelevant details. We conduct extensive experiments on two image-to-video translation tasks: facial expression retargeting and human pose forecasting. Superior results over the state-of-the-art methods on both tasks demonstrate the effectiveness of our approach.

Keywords: Video generation

•Motion forecasting •Residual learning

1

Teaching Machines to Understand

Baseball Games: Large-Scale Baseball

Video Database for Multiple Video

Understanding Tasks

Minho Shim , Young Hwi Kim , Kyungmin Kim , and Seon Joo Kim(B)

Yonsei University, Seoul, South Korea

{minhoshim,younghwikim,kyungminkim,seonjookim }@yonsei.ac.kr

Abstract. A major obstacle in teaching machines to understand videos is the lack of training data, as creating temporal annotations for long videos requires a huge amount of human effort. To this end, we introduce a new large-scale baseball video dataset called the BBDB, which is produced semi-automatically by using play-by-play texts available online. The BBDB contains 4200+hr of baseball game videos with 400k+ temporally annotated activity segments. The new dataset has several major challenging factors compared to other datasets: (1) the dataset contains a large number of visually similar segments with different labels. (2) It can be used for many video understanding tasks including video recognition, localization, text-video alignment, video highlight generation, and data imbalance problem. To observe the potential of the BBDB, we conducted extensive experiments by running many different types of video understanding algorithms on our new dataset. The database is available at <https://sites.google.com/site/eccv2018bbdb/>.

Keywords: Video understanding

•Large-scale video dataset

Action recognition •Temporal localization

1

SketchyScene: Richly-Annotated Scene

Sketches

Changqing Zou¹, Qian Yu², Ruofei Du¹, Haoran Mo³, Yi-Zhe Song²,

Tao Xiang², Chengying Gao³, Baoquan Chen^{4(B)}, and Hao Zhang⁵

¹University of Maryland, College Park, US

²Queen Mary University of London, London, UK

³Sun Yat-sen University, Guangzhou, China

4Shandong University, Jinan, China

baoquan@sdu.edu.cn

5Simon Fraser University, Burnaby, Canada

Abstract. We contribute the first large-scale dataset of scene sketches, SketchyScene, with the goal of advancing research on sketch understanding at both the object and scene level. The dataset is created through a novel and carefully designed crowdsourcing pipeline, enabling users to efficiently generate large quantities of realistic and diverse scene sketches. SketchyScene contains more than 29,000 scene-level sketches, 7,000+ pairs of scene templates and photos, and 11,000+ object sketches. All objects in the scene sketches have ground-truth semantic and instance masks. The dataset is also highly scalable and extensible, easily allowing augmenting and/or changing scene composition. We demonstrate the potential impact of SketchyScene by training new computational models for semantic segmentation of scene sketches and showing how the new dataset enables several applications including image retrieval, sketch colorization, editing, and captioning, etc. The dataset and code can be found at <https://github.com/SketchyScene/SketchyScene>.

Keywords: Sketch dataset

·Scene sketch ·Sketch segmentation

1

SketchyScene: Richly-Annotated Scene Sketches 441

TU-Berlin [9] is the first such large-scale crowd-sourced sketch dataset which was primarily designed for sketch recognition. It consists of 20,000 sketches spanning over 250 categories. The more recent QuickDraw [10] dataset is much larger,

with 50 million sketches across 345 categories. Albeit being large enough to facilitate

stroke-level analysis [6], sketches sourced in these datasets were produced by

sketching towards a semantic concept (e.g., "cat", "house"), without a reference photo or mental recollection of natural scene/objects. This greatly limits the level of visual detail and variations depicted, therefore making them unfitting for fine-grained matching and scene-level parsing. For example, faces are almost all in their frontal view, and depicted as a smiley in QuickDraw.

The concurrent work of [36] and [19] progressed the field further by collecting object instance sketches for FG-SBIR. QMUL database [36] consists of 716 sketch-photo pairs across two object categories (shoe and chair), with refer-

ence photos crawled from on-line shopping websites. Sketchy [19] contains 75,471

sketches and 12,500 corresponding photos across a much wider selection of categories (125 in total). Object instance sketches are produced by asking crowd-sourcers to depict their mental recollection of a reference photo. In comparison with concept sketches [9,10], they by and large exhibit more object details and have matching poses with the reference photos. However, a common drawback for both, for the purpose of this project, lies with their limited pose selection and object configurations. QMUL sketches exhibit only one object pose (side view) under a single object configuration. Scene sketches albeit exhibit more object poses and configurations, are still restricted since their reference photos mainly consist of single objects centered on relatively plain backgrounds (thus depict

no object interactions). This drawback essentially renders them both unsuitable for our task of scene sketch parsing, where complex mutual object interactions dictate high degree of object pose and configuration variations, as well as subtle

details. For example, within a picnic scene depicted in Figure 1, people appear in different poses and configurations with subtle eye contacts among each other. Fig.2(c) shows a composition result using sketches from Sketchy and TU-Berlin.

SketchyScene is the first large-scale dataset specifically designed for scene-level sketch understanding. It differs from all aforementioned datasets in that it goes beyond single object sketch understanding to tackle scene sketch, and purposefully includes an assorted selection of object sketches with diverse poses, configurations and object details to accommodate the complex scene-level object interactions. Although the existing dataset Abstract Scenes [38] serves a similar

motivation for understanding high-level semantic information in visual data, they focus on abstract scenes composed using clip arts, which include much more visual cues such as color and texture. In addition, their scenes are restricted in describing interactions between two characters and a handful of objects, while the scene contents and mutual object interactions in SketchyScene are a lot more diverse.

2.2 Sketch Understanding

Sketch recognition is perhaps the most studied problem in sketch understanding. Since the release of TU-Berlin dataset [9], many works have been proposed and

Learn-to-Score: Efficient 3D Scene

Exploration by Predicting View Utility

Benjamin Hepp^{1,2(B)}, Debadeepta Dey², Sudipta N. Sinha², Ashish Kapoor²,

Neel Joshi², and Otmar Hilliges¹

¹ETH Zurich, Zurich, Switzerland

benjamin.hepp@inf.ethz.ch

²Microsoft Research, Redmond, USA

Abstract. Camera equipped drones are nowadays being used to explore large scenes and reconstruct detailed 3D maps. When free space in the scene is approximately known, an offline planner can generate optimal plans to efficiently explore the scene. However, for exploring unknown scenes, the planner must predict and maximize usefulness of where to go on the fly. Traditionally, this has been achieved using handcrafted utility functions. We propose to learn a better utility function that predicts the usefulness of future viewpoints. Our learned utility function is based on a 3D convolutional neural network. This network takes as input a novel volumetric scene representation that implicitly captures previously visited viewpoints and generalizes to new scenes. We evaluate our method on several large 3D models of urban scenes using simulated depth cameras. We show that our method outperforms existing utility measures in terms of reconstruction performance and is robust to sensor noise.

Keywords: 3D reconstruction

· Exploration · Active vision · 3D CNN

1

Revisiting RCNN: On Awakening

the Classification Power of Faster RCNN

Bowen Cheng¹, Yunchao Wei^{1(B)}, Honghui Shi², Rogerio Feris²,

Jinjun Xiong², and Thomas Huang¹

¹University of Illinois at Urbana-Champaign, Champaign, IL, USA

{bcheng9,yunchao,t-huang1}@illinois.edu

²IBM T.J. Watson Research Center, Ossining, NY, USA

Honghui.Shi@ibm.com, {rsferis,jinjun}@us.ibm.com

Abstract. Recent region-based object detectors are usually built with separate classification and localization branches on top of shared feature extraction networks. In this paper, we analyze failure cases of state-of-the-art detectors and observe that most hard false positives result from classification instead of localization. We conjecture that: (1) Shared feature representation is not optimal due to the mismatched goals of feature learning for classification and localization; (2) multi-task learning helps, yet optimization of the multi-task loss may result in sub-optimal for individual tasks; (3) large receptive field for different scales leads to redundant context information for small objects. We demonstrate

the potential of detector classification power by a simple, effective, and widely-applicable Decoupled Classification Refinement (DCR) network. DCR samples hard false positives from the base classifier in Faster RCNN and trains a RCNN-styled strong classifier. Experiments show new state-of-the-art results on PASCAL VOC and COCO without any bells and whistles.

Keywords: Object detection

1

Semi-supervised Generative Adversarial

Hashing for Image Retrieval

Guan'an Wang^{1,3}, Qinghao Hu^{2,3}, Jian Cheng^{2,3,4},

and Zengguang Hou^{1,3,4}(B)

¹The State Key Laboratory for Management and Control of Complex Systems,
Institute of Automation, Chinese Academy of Sciences, Beijing, China

{wangguan2015,zengguang.hou}@ia.ac.cn

²National Laboratory of Pattern Recognition, Institute of Automation, Chinese
Academy of Sciences, Beijing, China

{qinghao.hu,jcheng}@nlpr.ia.ac.cn

³University of Chinese Academy of Sciences, Beijing, China

⁴Center for Excellence in Brain Science and Intelligence Technology, Beijing, China

Abstract. With explosive growth of image and video data on the Internet, hashing technique has been extensively studied for large-scale visual search. Benefiting from the advance of deep learning, deep hashing methods have achieved promising performance. However, those deep hashing models are usually trained with supervised information, which is rare and expensive in practice, especially class labels. In this paper, inspired by the idea of generative models and the minimax two-player game, we propose a novel semi-supervised generative adversarial hashing (SSGAH) approach. Firstly, we unify a generative model, a discriminative model and a deep hashing model in a framework for making use of triplet-wise information and unlabeled data. Secondly, we design novel structure of the generative model and the discriminative model to learn the distribution of triplet-wise information in a semi-supervised way. In addition, we propose a semi-supervised ranking loss and an adversary ranking loss to learn binary codes which preserve semantic similarity for both labeled data and unlabeled data. Finally, by optimizing the whole model in an adversary training way, the learned binary codes can capture better semantic information of all data. Extensive empirical evaluations on two widely-used benchmark datasets show that our proposed approach significantly outperforms state-of-the-art hashing methods.

Keywords: Information retrieval

· Hashing · Deep learning · GANs

1

Person Re-identification with Deep

Similarity-Guided Graph Neural Network

Yantao Shen¹, Hongsheng Li¹(B), Shuai Yi², Dapeng Chen¹,

and Xiaogang Wang¹

¹CUHK-SenseTime Joint Lab, The Chinese University of Hong Kong,
Hong Kong, People's Republic of China

{ytshen,hsli,dpchen}@ee.cuhk.edu.hk

²SenseTime Research, Hong Kong, People's Republic of China

yishuai@sensetime.com

Abstract. The person re-identification task requires to robustly estimate visual similarities between person images. However, existing person re-identification models mostly estimate the similarities of different image pairs of probe and gallery images independently while ignores the relationship information between different probe-gallery pairs. As a result, the similarity estimation of some hard samples might not be accurate. I

In this paper, we propose a novel deep learning framework, named Similarity-Guided Graph Neural Network (SGGNN) to overcome such limitations. Given a probe image and several gallery images, SGGNN creates a graph to represent the pairwise relationships between probe-gallery pairs (nodes) and utilizes such relationships to update the probe-gallery relation features in an end-to-end manner. Accurate similarity estimation can be achieved by using such updated probe-gallery relation features for prediction. The input features for nodes on the graph are the relation features of different probe-gallery image pairs. The probe-gallery relation feature updating is then performed by the messages passing in SGGNN, which takes other nodes' information into account for similarity estimation. Different from conventional GNN approaches, SGGNN learns the edge weights with rich labels of gallery instance pairs directly, which provides relation fusion more precise information. The effectiveness of our proposed method is validated on three public person re-identification datasets.

Keywords: Deep learning
 · Person re-identification
 Graph Neural Networks

1

Learning and Matching Multi-View
 Descriptors for Registration of Point
 Clouds

Lei Zhou¹, Siyu Zhang¹, Zixin Luo¹, Tianwei Shen¹, Runze Zhang^{1(B)},
 Mingmin Zhen¹, Tian Fang², and Long Quan¹

¹Hong Kong University of Science and Technology, Hong Kong, China
 {lzhouai, szhu, zluoag, tshenaa, rzhangaj, mzhen, quan}@cse.ust.hk

²Shenzhen Zhuke Innovation Technology (Altizure), Shenzhen, China
 fangtian@altizure.com

Abstract. Critical to the registration of point clouds is the establishment of a set of accurate correspondences between points in 3D space. The correspondence problem is generally addressed by the design of discriminative 3D local descriptors on the one hand, and the development of robust matching strategies on the other hand. In this work, we first propose a multi-view local descriptor, which is learned from the images of multiple views, for the description of 3D keypoints. Then, we develop a robust matching approach, aiming at rejecting outlier matches based on the efficient inference via belief propagation on the defined graphical model. We have demonstrated the boost of our approaches to registration on the public scanning and multi-view stereo datasets. The superior performance has been verified by the intensive comparisons against a variety of descriptors and matching methods.

Keywords: Point cloud registration
 · 3D descriptor · Robust matching

1

Revisiting Autofocus for Smartphone
 Cameras

Abdullah Abuolaim^(B), Abhijith Punnapurath, and Michael S. Brown
 Department of Electrical Engineering and Computer Science,
 Lassonde School of Engineering, York University, Toronto, Canada
 {abuolaim, pabhijith, mbrown}@eecs.yorku.ca

Abstract. Autofocus (AF) on smartphones is the process of determining how to move a camera's lens such that certain scene content is in focus. The underlying algorithms used by AF systems, such as contrast detection and phase differencing, are well established. However, determining a high-level objective regarding how to best focus a particular scene is less clear. This is evident in part by the fact that different smartphone cameras employ different AF criteria; for example, some attempt to keep item

s in the center in focus, others give priority to faces while others maximize the sharpness of the entire scene. The fact that different objectives exist raises the research question of whether there is a preferred objective. This becomes more interesting when AF is applied to videos of dynamic scenes. The work in this paper aims to revisit AF for smartphones within the context of temporal image data. As part of this effort, we describe the capture of a new 4D dataset that provides access to a full focal stack at each time point in a temporal sequence. Based on this dataset, we have developed a platform and associated application programming interface (API) that mimic real AF systems, restricting lens motion within the constraints of a dynamic environment and frame capture. Using our platform we evaluated several high-level focusing objectives and found interesting insight into what users prefer. We believe our new temporal focal stack dataset, AF platform, and initial user-study findings will be useful in advancing AF research.

Keywords: Autofocus

•Focal stack •AF platform

Low-level computer vision

1

Deep Burst Denoising

Clément Godard^{1(B)}, Kevin Matzen², and Matt Uyttendaele²

¹University College London, London, UK

c.godard@cs.ucl.ac.uk

²Facebook, Seattle, USA

Abstract. Noise is an inherent issue of low-light image capture, which is worsened on mobile devices due to their narrow apertures and small sensors. One strategy for mitigating noise in low-light situations is to increase the shutter time, allowing each photosite to integrate more light and decrease noise variance. However, there are two downsides of long exposures: (a) bright regions can exceed the sensor range, and (b) camera and scene motion will cause blur. Another way of gathering more light is to capture multiple short (thus noisy) frames in a burst and intelligently integrate the content, thus avoiding the above downsides. In this paper, we use the burst-capture strategy and implement the intelligent integration via a recurrent fully convolutional deep neural net (CNN). We build our novel, multi-frame architecture to be a simple addition to any single frame denoising model. The resulting architecture denoises all frames in a sequence of arbitrary length. We show that it achieves state of the art denoising results on our burst dataset, improving on the best published multi-frame techniques, such as VBM4D and FlexISP. Finally, we explore other applications of multi-frame image enhancement and show that our CNN architecture generalizes well to image super-resolution.

1

ISNN: Impact Sound Neural Network

for Audio-Visual Object Classification

Auston Sterling¹, Justin Wilson^{1(B)}, Sam Lowe¹, and Ming C. Lin^{1,2}

¹Department of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, USA

{austonst,wilson,samlowe,lin}@cs.unc.edu

²Department of Computer Science, University of Maryland, College Park, USA

lin@cs.umd.edu

Abstract. 3D object geometry reconstruction remains a challenge when working with transparent, occluded, or highly reflective surfaces. While recent methods classify shape features using raw audio, we present a multimodal neural network optimized for estimating an object's geometry and material. Our networks use spectrograms of recorded and synthesized object impact sounds and voxelized shape estimates to extend the capabilities of vision-based reconstruction. We evaluate our method on

multiple datasets of both recorded and synthesized sounds. We further present an interactive application for real-time scene reconstruction in which a user can strike objects, producing sound that can instantly classify and segment the struck object, even if the object is transparent or visually occluded.

1

StereoNet: Guided Hierarchical
Refinement for Real-Time Edge-Aware
Depth Prediction

Sameh Khamis(B), Sean Fanello, Christoph Rhemann, Adarsh Kowdle,
Julien Valentin, and Shahram Izadi
Google Inc., Mountain View, USA
sameh@google.com

Abstract. This paper presents StereoNet, the first end-to-end deep architecture for real-time stereo matching that runs at 60fps on an NVidia Titan X, producing high-quality, edge-preserved, quantization-free disparity maps. A key insight of this paper is that the network achieves a sub-pixel matching precision that is a magnitude higher than those of traditional stereo matching approaches. This allows us to achieve real-time performance by using a very low resolution cost volume that encodes all the information needed to achieve high disparity precision. Spatial precision is achieved by employing a learned edge-aware upsampling function. Our model uses a Siamese network to extract features from the left and right image. A first estimate of the disparity is computed in a very low resolution cost volume, then hierarchically the model introduces high-frequency details through a learned upsampling function that uses compact pixel-to-pixel refinement networks. Leveraging color input as a guide, this function is capable of producing high-quality edge-aware output. We achieve compelling results on multiple benchmarks, showing how the proposed method offers extreme flexibility at an acceptable computational budget.

Keywords: Stereo matching

•Depth estimation

Edge-aware refinement •Cost volume filtering •Deep learning

1

Attention-Aware Deep Adversarial
Hashing for Cross-Modal Retrieval

Xi Zhang^{1,2}, Hanjiang Lai^{1,2(B)}, and Jiashi Feng³

¹School of Data and Computer Science, Sun Yat-Sen University, Guangzhou, China
zhangx368@mail2.sysu.edu.cn,

laihanj3@mail.sysu.edu.cn

²Guangdong Key Laboratory of Big Data Analysis and Processing,
Guangzhou, China

³Department of Electrical and Computer Engineering, National University
of Singapore, Singapore, Singapore

elefjia@nus.edu.sg

Abstract. Due to the rapid growth of multi-modal data, hashing methods for cross-modal retrieval have received considerable attention. However, finding content similarities between different modalities of data is still challenging due to an existing heterogeneity gap. To further address this problem, we propose an adversarial hashing network with an attention mechanism to enhance the measurement of content similarities by selectively focusing on the informative parts of multi-modal data. The proposed new deep adversarial network consists of three building blocks: (1) the feature learning module to obtain the feature representations; (2) the attention module to generate an attention mask, which is used to divide the feature representations into the attended and unattended feature representations; and (3) the hashing module to learn hash functions that preserve the

e similarities between different modalities. In our framework, the attention and hashing modules are trained in an adversarial way: the attention module attempts to make the hashing module unable to pre-serve the similarities of multi-modal data w.r.t. the unattended feature representations, while the hashing module aims to preserve the similarities of multi-modal data w.r.t. the attended and unattended feature representations. Extensive evaluations on several benchmark datasets demonstrate that the proposed method brings substantial improvements over other state-of-the-art cross-modal hashing methods.

Keywords: Hashing

·Adversarial learning ·Attention mechanism

Cross modal retrieval

1

3DFeat-Net: Weakly Supervised Local 3D

Features for Point Cloud Registration

Zi Jian Yew(B)and Gim Hee Lee

Department of Computer Science, National University of Singapore,
Singapore, Singapore

{zijian.yew,gimhee.lee}@comp.nus.edu.sg

Abstract. In this paper, we propose the 3DFeat-Net which learns both 3D feature detector and descriptor for point cloud matching using weak supervision. Unlike many existing works, we do not require manual annotation of matching point clusters. Instead, we leverage on alignment and attention mechanisms to learn feature correspondences from GPS/INS tagged 3D point clouds without explicitly specifying them. We create training and benchmark outdoor Lidar datasets, and experiments show that 3DFeat-Net obtains state-of-the-art performance on these gravity-aligned datasets.

Keywords: Point cloud

·Registration ·Deep learning

Weak supervision

1

Deep Domain Generalization via
Conditional Invariant Adversarial
Networks

Ya Li1, Xinmei Tian1(B), Mingming Gong2,3, Yajing Liu1, Tongliang Liu4,
Kun Zhang2, and Dacheng Tao4

1CAS Key Laboratory of Technology in Geo-Spatial Information
Processing and Application Systems,

University of Science and Technology of China, Hefei, China

{muziye,lyj123}@mail.ustc.edu.cn, xinmei@ustc.edu.cn

2Department of Philosophy, Carnegie Mellon University, Pittsburgh, USA
gongmingnju@gmail.com, kunz1@cmu.edu

3Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh, USA

4UBTECH Sydney AI Centre, SIT, FEIT, University of Sydney, Sydney, Australia
tliang.liu@gmail.com, dacheng.tao@sydney.edu.au

Abstract. Domain generalization aims to learn a classification model from multiple source domains and generalize it to unseen target domains. A critical problem in domain generalization involves learning domain-invariant representations. Let X and Y denote the features and the labels, respectively. Under the assumption that the conditional distribution $P(Y|X)$ remains unchanged across domains, earlier approaches to domain generalization learned the invariant representation $T(X)$ by minimizing the discrepancy of the marginal distribution $P(T(X))$. However, such an assumption of stable $P(Y|X)$ does not necessarily hold in practice. In addition, the representation learning function $T(X)$ is usually constrained to a simple linear transformation or shallow networks. To address the above two drawbacks, we propose an end-to-end conditional

invariant deep domain generalization approach by leveraging deep neural networks for domain-invariant representation learning. The domain-invariance property is guaranteed through a conditional invariant adversarial network that can learn domain-invariant representations w.r.t. the joint distribution $P(T(X), Y)$ if the target domain data are not severely class unbalanced. We perform various experiments to demonstrate the effectiveness of the proposed method.

Keywords: Domain generalization

•Adversarial networks

Domain invariant representation

c/circlecopyrtSpringer Nature Switzerland AG 2018

V. Ferrari et al. (Eds.): ECCV 2018, LNCS 11219, pp. 647–663, 2018.<https://doi.org/10.1007/978-3-030-01267-0>

_3

Using LIP to Gloss Over Faces

in Single-Stage Face Detection Networks

Siqi Yang(B), Arnold Wiliem, Shaokang Chen, and Brian C. Lovell

The University of Queensland, Brisbane, Australia

siqi.yang@uq.net.au , a.wiliem@uq.edu.au , shaokangchenuq@gmail.com ,

lovell@itee.uq.edu.au

Abstract. This work shows that it is possible to fool/attack recent state-of-the-art face detectors which are based on the single-stage networks. Successfully attacking face detectors could be a serious malware vulnerability when deploying a smart surveillance system utilizing face detectors. In addition, for the privacy concern, it helps prevent faces being harvested and stored in the server. We show that existing adversarial perturbation methods are not effective to perform such an attack, especially when there are multiple faces in the input image. This is because the adversarial perturbation specifically generated for one face may disrupt the adversarial perturbation for another face. In this paper, we call this problem the Instance Perturbation Interference (IPI) problem. This IPI problem is addressed by studying the relationship between the deep neural network receptive field and the adversarial perturbation. Besides the single-stage face detector, we find that the IPI problem also exists on the first stage of the Faster-RCNN, the commonly used two-stage object detector. As such, we propose the Localized Instance Perturbation (LIP) that confines the adversarial perturbation inside the Effective Receptive Field (ERF) of a target to perform the attack. Experimental results show the LIP method massively outperforms existing adversarial perturbation generation methods – often by a factor of 2 to 10.

Keywords: Adversarial

•Interference •Effective Receptive Field

Single-stage network •Detection

1

HiDDeN: Hiding Data With Deep

Networks

Jiren Zhu(B), Russell Kaplan, Justin Johnson, and Li Fei-Fei

Computer Science Department, Stanford University, Stanford, USA

{jirenz, rjkaplan, jcjohns, feifeili}@cs.stanford.edu

Abstract. Recent work has shown that deep neural networks are highly sensitive to tiny perturbations of input images, giving rise to adversarial examples. Though this property is usually considered a weakness of learned models, we explore whether it can be beneficial. We find that neural networks can learn to use invisible perturbations to encode a rich amount of useful information. In fact, one can exploit this capability for the task of data hiding. We jointly train encoder and decoder networks, where given an input message and cover image, the encoder produces a visually indistinguishable encoded image, from which the decoder can

recover the original message. We show that these encodings are competitive with existing data hiding algorithms, and further that they can be made robust to noise: our models learn to reconstruct hidden information in an encoded image despite the presence of Gaussian blurring, pixel-wise dropout, cropping, and JPEG compression. Even though JPEG is non-differentiable, we show that a robust model can be trained using differentiable approximations. Finally, we demonstrate that adversarial training improves the visual quality of encoded images.

Keywords: Adversarial networks

·Steganography

Robust blind watermarking ·Deep learning ·Convolutional networks

1

Multimodal Dual Attention Memory for

Video Story Question Answering

Kyung-Min Kim¹, Seong-Ho Choi², Jin-Hwa Kim³,

and Byoung-Tak Zhang^{2,4(B)}

¹Clova AI Research, NAVER Corp, Seongnam 13561, South Korea

kmkim@bi.snu.ac.kr

²Computer Science and Engineering, Seoul National University,

Seoul 08826, South Korea

{shchoi,btzhang}@bi.snu.ac.kr

³SK T-Brain, Seoul 04539, South Korea

jhkim@bi.snu.ac.kr

⁴Surromind Robotics, Seoul 08826, South Korea

Abstract. We propose a video story question-answering (QA) architecture, Multimodal Dual Attention Memory (MDAM). The key idea is to use a dual attention mechanism with late fusion. MDAM uses self-attention to learn the latent concepts in scene frames and captions. Given a question, MDAM uses the second attention over these latent concepts. Multimodal fusion is performed after the dual attention processes (late fusion). Using this processing pipeline, MDAM learns to infer a high-level vision-language joint representation from an abstraction of the full video content. We evaluate MDAM on PororoQA and MovieQA datasets which have large-scale QA annotations on cartoon videos and movies, respectively. For both datasets, MDAM achieves new state-of-the-art results with significant margins compared to the runner-up models. We confirm the best performance of the dual attention mechanism combined with late fusion by ablation studies. We also perform qualitative analysis by visualizing the inference mechanisms of MDAM.

Keywords: Video story QA

·Visual QA ·Attention mechanism

Multimodal learning ·Deep learning

1

Deep Variational Metric Learning

Xudong Lin, Yueqi Duan, Qiyuan Dong, Jiwen Lu(B), and Jie Zhou

Tsinghua University, Beijing 100084, China

stein6@163.com, {duanyq14,dqy15}@mails.tsinghua.edu.cn,

{lujiwen,jzhou}@tsinghua.edu.cn

Abstract. Deep metric learning has been extensively explored recently, which trains a deep neural network to produce discriminative embedding features. Most existing methods usually enforce the model to be indiscriminating to intra-class variance, which makes the model overfitting to the training set to minimize loss functions on these specific changes and leads to low generalization power on unseen classes. However, these methods ignore a fact that in the central latent space, the distribution of variance within classes is actually independent on classes. In this paper, we propose a deep variational metric learning (DVML) framework to explicitly model the intra-class variance and disentangle the intra-class

invariance, namely, the class centers. With the learned distribution of intra-class variance, we can simultaneously generate discriminative samples to improve robustness. Our method is applicable to most of existing metric learning algorithms, and extensive experiments on three benchmark datasets including CUB-200-2011, Cars196 and Stanford Online Products show that our DVML significantly boosts the performance of currently popular deep metric learning methods.

Keywords: Metric learning

• Variational auto-encoder

Discriminative samples generating

1

HGMR: Hierarchical Gaussian Mixtures

for Adaptive 3D Registration

Benjamin Eckart(B), Kihwan Kim, and Jan Kautz

NVIDIA Research, Santa Clara, CA, USA

beckart@nvidia.com

Abstract. Point cloud registration sits at the core of many important and challenging 3D perception problems including autonomous navigation, SLAM, object/scene recognition, and augmented reality. In this paper, we present a new registration algorithm that is able to achieve state-of-the-art speed and accuracy through its use of a Hierarchical Gaussian Mixture representation. Our method, Hierarchical Gaussian Mixture Registration (HGMR), constructs a top-down multi-scale representation of point cloud data by recursively running many small-scale data likelihood segmentations in parallel on a GPU. We leverage the resulting representation using a novel optimization criterion that adaptively finds the best scale to perform data association between spatial subsets of point cloud data. Compared to previous Iterative Closest Point and GMM-based techniques, our tree-based point association algorithm performs data association in logarithmic-time while dynamically adjusting the level of detail to best match the complexity and spatial distribution characteristics of local scene geometry. In addition, unlike other GMM methods that restrict covariances to be isotropic, our new PCA-based optimization criterion well-approximates the true MLE solution even when fully anisotropic Gaussian covariances are used. Efficient data association, multi-scale adaptability, and a robust MLE approximation produce an algorithm that is up to an order of magnitude both faster and more accurate than current state-of-the-art on a wide variety of 3D datasets captured from LiDAR to structured light.

1

Bi-Real Net: Enhancing the Performance

of 1-Bit CNNs with Improved

Representational Capability

and Advanced Training Algorithm

Zechun Li¹, Baoyuan Wu², Wenhan Luo², Xin Yang³(B), Wei Liu²,

and Kwang-Ting Cheng¹

¹Hong Kong University of Science and Technology, Hong Kong, China

zliubq@connect.ust.hk, timcheng@ust.hk

²Tencent AI Lab, Beijing, China

wubaoyuan1987@gmail.com, wluo.china@gmail.com, wliu@ee.columbia.edu

³Huazhong University of Science and Technology, Wuhan, China

xinyang2014@hust.edu.cn

Abstract. In this work, we study the 1-bit convolutional neural networks (CNNs), of which both the weights and activations are binary. While being efficient, the classification accuracy of the current 1-bit CNNs is much worse compared to their counterpart real-valued CNN models on the large-scale dataset, like ImageNet. To minimize the performance gap between the 1-bit and real-valued CNN models, we propose

a novel model, dubbed Bi-Real net, which connects the real activations (after the 1-bit convolution and/or BatchNorm layer, before the signfunction) to activations of the consecutive block, through an identity shortcut. Consequently, compared to the standard 1-bit CNN, the representational capability of the Bi-Real net is significantly enhanced and the additional cost on computation is negligible. Moreover, we develop a specific training algorithm including three technical novelties for 1-bit CNNs. Firstly, we derive a tight approximation to the derivative of the non-differentiable sign function with respect to activation. Secondly, we propose a magnitude-aware gradient with respect to the weight for updating the weight parameters. Thirdly, we pre-train the real-valued CNN model with a clip function, rather than the ReLU function, to better initialize the Bi-Real net. Experiments on ImageNet show that the Bi-Real net with the proposed training algorithm achieves 56.4% and 62.2% top-1 accuracy with 18 layers and 34 layers, respectively. Compared to the state-of-the-arts (e.g., XNOR Net), Bi-Real net achieves up to 10% higher top-1 accuracy with more memory saving and lower computational cost.

©/circlecopyrt Springer Nature Switzerland AG 2018

V. Ferrari et al. (Eds.): ECCV 2018, LNCS 11219, pp. 747–763, 2018. <https://doi.org/10.1007/978-3-030-01267-0>

_4

Orthogonal Deep Features Decomposition

for Age-Invariant Face Recognition

Yitong Wang, Dihong Gong, Zheng Zhou, Xin GJi, Haowang,

Zhifeng Li(B), Wei Liu(B), and Tong Zhang

Tencent AI Lab, Beijing, China

{yitongwang, encorezhou, denisji, hawelwang, michaelzfli}@tencent.com,

gongdihong@gmail.com, wl2223@columbia.edu, tongzhang@tongzhang-ml.org

Abstract. As facial appearance is subject to significant intra-class variations caused by the aging process over time, age-invariant face recognition (AIFR) remains a major challenge in face recognition community. To reduce the intra-class discrepancy caused by the aging, in this paper we propose a novel approach (namely, Orthogonal Embedding CNNs, or OE-CNNs) to learn the age-invariant deep face features. Specifically, we decompose deep face features into two orthogonal components to represent age-related and identity-related features. As a result, identity-related features that are robust to aging are then used for AIFR. Besides, for complementing the existing cross-age datasets and advancing the research in this field, we construct a brand-new large-scale Cross-Age Face dataset (CAF). Extensive experiments conducted on the three public domain face aging datasets (MORPH Album 2, CACD-VS and FG-NET) have shown the effectiveness of the proposed approach and the value of the constructed CAF dataset on AIFR. Benchmarking our algorithm on one of the most popular general face recognition (GFR) dataset LFW additionally demonstrates the comparable generalization performance on GFR.

Keywords: Age-invariant face recognition

Convolutional neural networks

• Cross-age face dataset

1

Broadcasting Convolutional Network for

Visual Relational Reasoning

Simyung Chang^{1,2}, John Yang¹, SeongUk Park¹, and Nojun Kwak^{1(B)}

¹Seoul National University, Seoul, South Korea

{yjohn, swpark0703, nojunk}@snu.ac.kr

²Samsung Electronics, Suwon, South Korea

timelighter@snu.ac.kr

Abstract. In this paper, we propose the Broadcasting Convolutional

Network (BCN) that extracts key object features from the global field of an entire input image and recognizes their relationship with local features. BCN is a simple network module that collects effective spatial features, embeds location information and broadcasts them to the entire feature maps. We further introduce the Multi-Relational Network (multiRN) that improves the existing Relation Network (RN) by utilizing the BCN module. In pixel-based relation reasoning problems, with the help of BCN, multiRN extends the concept of 'pairwise relations' in conventional RNs to 'multiwise relations' by relating each object with multiple objects at once. This yields in $O(n)$ complexity for n objects, which is a vast computational gain from RNs that take $O(n^2)$. Through

experiments, multiRN has achieved a state-of-the-art performance on CLEVR dataset, which proves the usability of BCN on relation reasoning problems.

Keywords: Visual relational reasoning · BCN · Broadcast · CLEVR

Multi-RN · Visuo-spatial features

1

Improving Spatiotemporal

Self-supervision by Deep Reinforcement

Learning

Uta Buechler(B), Biagio Brattoli(B), and Björn Ommer

Heidelberg University, HCI/IWR, Heidelberg, Germany

{uta.buechler,biagio.brattoli,bjoern.ommer}@iwr.uni-heidelberg.de

Abstract. Self-supervised learning of convolutional neural networks can harness large amounts of cheap unlabeled data to train powerful feature representations. As surrogate task, we jointly address ordering of visual data in the spatial and temporal domain. The permutations of training samples, which are at the core of self-supervision by ordering, have so far been sampled randomly from a fixed preselected set. Based on deep reinforcement learning we propose a sampling policy that adapts to the state of the network, which is being trained. Therefore, new permutations are sampled according to their expected utility for updating the convolutional feature representation. Experimental evaluation on unsupervised and transfer learning tasks demonstrates competitive performance on standard benchmarks for image and video classification and nearest neighbor retrieval.

Keywords: Deep reinforcement learning

· Self-supervision · Shuffling

Action recognition · Image understanding

1

Learning to Look around Objects

for Top-View Representations of Outdoor

Scenes

Samuel Schulter¹(B), Menghua Zhai², Nathan Jacobs²,

and Manmohan Chandraker^{1,3}

¹NEC-Laboratories, Cupertino, CA 95014, USA

samuel@nec-labs.com

²University of Kentucky, Lexington, KY 40506, USA

³University of California San Diego, La Jolla, CA 92093, USA

Abstract. Given a single RGB image of a complex outdoor road scene in the perspective view, we address the novel problem of estimating an occlusion-reasoned semantic scene layout in the top-view. This challenging problem not only requires an accurate understanding of both the 3D geometry and the semantics of the visible scene, but also of occluded areas. We propose a convolutional neural network that learns to predict occluded portions of the scene layout by looking around foreground objects like cars or pedestrians. But instead of hallucinating RGB val-

ues, we show that directly predicting the semantics and depths in the occluded areas enables a better transformation into the top-view. We further show that this initial top-view representation can be significantly enhanced by learning priors and rules about typical road layouts from simulated or, if available, map data. Crucially, training our model does not require costly or subjective human annotations for occluded areas or the top-view, but rather uses readily available annotations for standard semantic segmentation in the perspective view. We extensively evaluate and analyze our approach on the KITTI and Cityscapes data sets.

Keywords: 3D scene understanding

·Occlusion reasoning

Semantic top-view representations

1

Hierarchical Metric Learning
and Matching for 2D and 3D Geometric
Correspondences

Mohammed E. Fathy¹, Quoc-Huy Tran^{2(B)}, M. Zeeshan Zia³, Paul Vernaza²,
and Manmohan Chandraker^{2,4}

¹Google Cloud AI, New York, USA

²NEC Laboratories America, Inc., Princeton, USA

qhtran@nec.labs.com

³Microsoft HoloLens, London, UK

⁴University of California, San Diego, US

Abstract. Interest point descriptors have fueled progress on almost every problem in computer vision. Recent advances in deep neural networks have enabled task-specific learned descriptors that outperform hand-crafted descriptors on many problems. We demonstrate that commonly used metric learning approaches do not optimally leverage the feature hierarchies learned in a Convolutional Neural Network (CNN), especially when applied to the task of geometric feature matching. While a metric loss applied to the deepest layer of a CNN, is often expected to yield ideal features irrespective of the task, in fact the growing receptive field as well as striding effects cause shallower features to be better at high precision matching tasks. We leverage this insight together with explicit supervision at multiple levels of the feature hierarchy for better regularization, to learn more effective descriptors in the context of geometric matching tasks. Further, we propose to use activation maps at different layers of a CNN, as an effective and principled replacement for the multi-resolution image pyramids often used for matching tasks. We propose concrete CNN architectures employing these ideas, and evaluate them on multiple datasets for 2D and 3D geometric matching as well as optical flow, demonstrating state-of-the-art results and generalization across datasets.

Keywords: Hierarchical metric learning

·Hierarchical matching

Geometric correspondences ·Dense correspondences

Part of this work was done during M. E. Fathy's internship at NEC Labs America. Code and models will be made available at <http://www.nec-labs.com/~mas/HiLM/>.

Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-01267-0_49) contains supplementary material, which is available to authorized users.

© Springer Nature Switzerland AG 2018

V. Ferrari et al. (Eds.): ECCV 2018, LNCS 11219, pp. 832–850, 2018. https://doi.org/10.1007/978-3-030-01267-0_49

_4

Deep Component Analysis

via Alternating Direction Neural

Networks

Calvin Murdock(B), Ming-Fang Chang, and Simon Lucey
Carnegie Mellon University, Pittsburgh, USA
{cmurdock,mingfanc,slucey}@cs.cmu.edu

Abstract. Despite a lack of theoretical understanding, deep neural networks have achieved unparalleled performance in a wide range of applications. On the other hand, shallow representation learning with component analysis is associated with rich intuition and theory, but smaller capacity often limits its usefulness. To bridge this gap, we introduce Deep Component Analysis (DeepCA), an expressive multilayer model formulation that enforces hierarchical structure through constraints on latent variables in each layer. For inference, we propose a differentiable optimization algorithm implemented using recurrent Alternating Direction Neural Networks (ADNNs) that enable parameter learning using standard backpropagation. By interpreting feed-forward networks as single-iteration approximations of inference in our model, we provide both a novel perspective for understanding them and a practical technique for constraining predictions with prior knowledge. Experimentally, we demonstrate performance improvements on a variety of tasks, including single-image depth prediction with sparse output constraints.

Keywords: Component analysis

·Deep learning ·Constraints

1

ADVISE: Symbolism and External
Knowledge for Decoding Advertisements

Keren Ye(B) and Adriana Kovashka
University of Pittsburgh, Pittsburgh, PA 15260, USA
{yekeren,kovashka}@cs.pitt.edu

Abstract. In order to convey the most content in their limited space, advertisements embed references to outside knowledge via symbolism. For example, a motorcycle stands for adventure (a positive property the ad wants associated with the product being sold), and a gun stands for danger (a negative property to dissuade viewers from undesirable behaviors). We show how to use symbolic references to better understand the meaning of an ad. We further show how anchoring ad understanding in general-purpose object recognition and image captioning improves results. We formulate the ad understanding task as matching the ad image to human-generated statements that describe the action that the ad prompts, and the rationale it provides for taking this action. Our proposed method outperforms the state of the art on this task, and on an alternative formulation of question-answering on ads. We show additional applications of our learned representations for matching ads to slogans, and clustering ads according to their topic, without extra training.

Keywords: Advertisements

·Symbolism ·Question answering

External knowledge ·Vision and language ·Representation learning

1
