

## XXIV Table of Contents

Joint Unsupervised Face Alignment and Behaviour Analysis .....	167
Lazaros Zafeiriou, Epameinondas Antonakos, Stefanos Zafeiriou, and Maja Pantic	
Learning a Deep Convolutional Network for Image Super-Resolution ....	184
Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang	
Discriminative Indexing for Probabilistic Image Patch Priors .....	200
Yan Wang, Sunghyun Cho, Jue Wang, and Shih-Fu Chang	
Modeling Video Dynamics with Deep Dynencoder .....	215
Xing Yan, Hong Chang, Shiguang Shan, and Xilin Chen	
Good Image Priors for Non-blind Deconvolution: Generic vs. Specific ...	231
Libin Sun, Sunghyun Cho, Jue Wang, and James Hays	
Image Deconvolution Ringing Artifact Detection and Removal via PSF	
Frequency Analysis .....	247
Ali Mosleh, J.M. Pierre Langlois, and Paul Green	
View-Consistent 3D Scene Flow Estimation over Multiple Frames .....	263
Christoph Vogel, Stefan Roth, and Konrad Schindler	
Hand Waving Away Scale .....	279
Christopher Ham, Simon Lucey, and Surya Singh	
A Non-Linear Filter for Gyroscope-Based Video Stabilization .....	294
Steven Bell, Alejandro Troccoli, and Kari Pulli	
Multi-modal and Multi-spectral Registration for Natural Images .....	309
Xiaoyong Shen, Li Xu, Qi Zhang, and Jiaya Jia	
Using Isometry to Classify Correct/Incorrect 3D-2D Correspondences ...	325
Toby Collins and Adrien Bartoli	
Bilateral Functions for Global Motion Modeling .....	341
Wen-Yan Daniel Lin, Ming-Ming Cheng, Jiangbo Lu,	
Hongsheng Yang, Minh N. Do, and Philip Torr	
VCDB: A Large-Scale Database for Partial Copy Detection in Videos ...	357
Yu-Gang Jiang, Yudong Jiang, and Jiajun Wang	
Single-Image Super-Resolution: A Benchmark .....	372
Chih-Yuan Yang, Chao Ma, and Ming-Hsuan Yang	
Well Begun Is Half Done: Generating High-Quality Seeds for Automatic	
Image Dataset Construction from Web .....	387
Yan Xia, Xudong Cao, Fang Wen, and Jian Sun	
Zero-Shot Learning via Visual Abstraction .....	401
Stanislaw Antol, C. Lawrence Zitnick, and Devi Parik	

\*\*\*\*\*

### Schwarps: Locally Projective Image Warps

Based on 2D Schwarzian Derivatives

Rahat Khan, Daniel Pizarro, and Adrien Bartoli

ISIT, UMR 6284 CNRS-UdA, Clermont-Ferrand, France

Abstract. Image warps -or just warps- capture the geometric deformation existing between two images of a deforming surface. The current approach to enforce a warp's smoothness is to penalize its second order partial derivatives. Because this favors locally affine warps, this fails to capture the local projective component of the image deformation. This may have a negative impact on applications such as image registration and deformable 3D reconstruction. We propose a novel penalty designed to smooth the warp while capturing the deformation's local projective structure. Our penalty is based on equivalents to the Schwarzian derivatives, which are projective differential invariants exactly preserved by homographies. We propose a methodology to derive a set of Partial Differential Equations with homographies as solutions. We call this system the Schwarzian equations and we explicitly derive them for 2D functions using differential properties of homographies. We name as Schwarp a warp which is estimated by penalizing the residual of Schwarzian equations. Experimental evaluation shows that Schwarps outperform existing warps in modeling and extrapolation power, and lead to far better results in Shape-from-Template and camera calibration from a deformable surface.

Keywords: Schwarzian Penalizer, Bending Energy, Projective Differential Invariants, Image Warps.

1

\*\*\*\*\*

gDLS: A Scalable Solution

to the Generalized Pose and Scale Problem

Chris Sweeney, Victor Fragoso, Tobias Höllerer, and Matthew Turk

University of California, Santa Barbara, USA

{cmsweeney,vfragoso,hollerer,mturk}@cs.ucsb.edu

Abstract. In this work, we present a scalable least-squares solution for computing a seven degree-of-freedom similarity transform. Our method utilizes the generalized camera model to compute relative rotation, translation, and scale from four or more 2D-3D correspondences. In particular, structure and motion estimations from monocular cameras lack scale without specific calibration. As such, our methods have applications in loop closure in visual odometry and registering multiple structure from motion reconstructions where scale must be recovered. We formulate the generalized pose and scale problem as a minimization of a least squares cost function and solve this minimization without iterations or initialization. Additionally, we obtain all minima of the cost function. The order of the polynomial system that we solve is independent of the number of points, allowing our overall approach to scale favorably. We evaluate our method experimentally on synthetic and real datasets and demonstrate that our methods produce higher accuracy similarity transform solutions than existing methods.

1

\*\*\*\*\*

Generalized Connectivity Constraints

for Spatio-temporal 3D Reconstruction

Martin Ralf Oswald, Jan Stuhmer, and Daniel Cremers

Department of Computer Science, Technische Universität München

Boltzmannstr. 3, 85748 Garching, Germany

Abstract. This paper introduces connectivity preserving constraints into spatio-temporal multi-view reconstruction. We efficiently model connectivity constraints by precomputing a geodesic shortest path tree on the occupancy likelihood. Connectivity of the final occupancy labeling is ensured with a set of linear constraints on the labeling function. In order to generalize the connectivity constraints from objects with genus 0 to an arbitrary genus, we detect loops by analyzing the visual hull of the scene. A modification of the constraints ensures connectivity in the presence of loops. The proposed efficient implementation adds little runtime and memory overhead to the reconstruction method. Several experiments show significant improvement over state-of-the-art methods and validate the practical use of this approach in scenes with fine structured details.

Keywords: connectivity constraints, spatio-temporal 3D reconstruction.

1 of 16 input images No Connectivity With a Connectivity Generalized Connectivity Constraint [22] Constraint [25]+[22] tivity Constraint

Fig. 1. Embedding connectivity constraints into multi-view reconstruction clearly helps to recover fine structures like the rope. The tree-shaped connectivity prior [25] only

works for objects without holes (genus 0), resulting in disconnected parts when the

rope touches the head. The proposed generalized connectivity constraint works for

objects with arbitrary genus. Dataset: 'jumping rope' sequence from the INRIA 4D repository [16].

This work was supported by the ERC Starting Grant 'Convex Vision' and the Technische Universität München - Institute for Advanced Study, funded by the German

Excellence Initiative.

D. Fleet et al. (Eds.): ECCV 2014, Part IV, LNCS 8692, pp. 32–46, 2014.  
c/circlecopyrtSpringer International Publishing Switzerland 201

\*\*\*\*\*

Passive Tomography of Turbulence Strength

Marina Alterman<sup>1</sup>, Yoav Y. Schechner<sup>1</sup>,  
Minh Vo<sup>2</sup>, and Srinivasa G. Narasimhan<sup>2</sup>

<sup>1</sup>Dept. Electrical Eng., Technion – Israel Institute of Technology, Haifa, Israel

<sup>2</sup>Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA

**Abstract.** Turbulence is studied extensively in remote sensing, astronomy, meteorology, aerodynamics and fluid dynamics. The strength of turbulence is a statistical measure of local variations in the turbulent medium. It influences engineering decisions made in these domains. Turbulence strength (TS) also affects safety of aircraft and tethered balloons, and reliability of free-space electromagnetic relays. We show that it is possible to estimate TS, without having to reconstruct instantaneous fluid flows. Instead, the TS field can be directly recovered, passively, using videos captured from different viewpoints. We formulate this as a linear tomography problem with a structure unique to turbulence fields. No tight synchronization between cameras is needed. Thus, realization is very simple to deploy using consumer-grade cameras. We experimentally demonstrate this both in a lab and in a large-scale uncontrolled complex outdoor environment, which includes industrial, rural and urban areas.

**1 The Need to Recover Turbulence Strength**

Turbulence creates refractive perturbations to light passing through a scene. This causes random distortions when imaging background objects. Hence, modeling and trying to compensate for random refractive distortions has long been studied in remote sensing [40], astronomy [34] and increasingly in computer vision [2,4,10,14,18,35,38,41,52,55]. Nevertheless, these distortions are not necessarily a problem: they offer information about the medium and the scene itself [44]. This insight is analogous to imaging in scattering media (fog [29], haze [19,37], water [11,30]), where visibility reduction yields ranging and characterizing of the

medium. Similar efforts are made to reconstruct refracting (transparent) solids or water surfaces [3,16,28,43,46] from images of a distorted background or light field [50,51]. In turbulence, refraction occurs continuously throughout a volume. We exploit random image distortions as a means to estimate the spatial (volumetric) distribution of turbulence strength (TS). The strength of turbulence is a statistical measure of local variations in the medium [20,21]. Often, it is not necessary to estimate an instantaneous snapshot of air density or refraction field [32,42]. Rather local statistics is relied upon heavily in many applications.

Meteorologists rely on TS to understand convection (which forms clouds), wind, and atmospheric stability. This is measured using special Doppler lidars [9,31], which are very expensive. Turbulence significantly affects the efficiency of wind

D. Fleet et al. (Eds.): ECCV 2014, Part IV, LNCS 8692, pp. 47–60, 2014.

c/circlecopyrtSpringer International Publishing Switzerland 201

\*\*\*\*\*

A Non-local Method

for Robust Noisy Image Completion

Wei Li, Lei Zhao, Duanqing Xu, and Dongming Lu  
Zhejiang University, Hangzhou, China

**Abstract.** The problem of noisy image completion refers to recovering an image from a random subset of its noisy intensities. In this paper, we propose a non-local patch-based algorithm to settle the noisy image completion problem following the methodology “grouping and collaboratively filtering”. The target of “grouping” is to form patch matrices by matching and stacking similar image patches. And the “collaboratively filtering” is achieved by transforming the tasks of simultaneously estimating missing values and removing noises for the stacked patch matrices.

ces into low-rank matrix completion problems, which can be efficiently solved by minimizing the nuclear norm of the matrix with linear constraints. The final output is produced by synthesizing all the restored patches. To improve the robustness of our algorithm, we employ an efficient and accurate patch matching method with adaptations including pre-completion and outliers removal, etc. Experiments demonstrate that our approach achieves state-of-the-art performance for the noisy image completion problem in terms of both PSNR and subjective visual quality.

1

\*\*\*\*\*

Improved Motion Invariant Deblurring  
through Motion Estimation

Scott McCloskey

Honeywell Labs, USA

Abstract. We address the capture of sharp images of fast-moving objects, and build on the Motion Invariant photographic technique. The key advantage of motion invariance is that, unlike other computational photographic techniques, it does not require pre-exposure velocity estimation in order to ensure numerically stable deblurring. Its disadvantage is that the invariance is only approximate - objects moving with non-zero velocity will exhibit artifacts in the deblurred image related to tail clipping in the motion Point Spread Function (PSF). We model these artifacts as a convolution of the desired latent image with an error PSF, and demonstrate that the spatial scale of these artifacts corresponds to the object velocity. Surprisingly, despite the use of parabolic motion to capture an image in which blur is invariant to motion, we demonstrate that the motion invariant image can be used to estimate object motion post-capture. With real camera images, we demonstrate significant reductions in the artifacts by using the estimated motion for deblurring. We also quantify a 96% reduction in reconstruction error, relative to a floor established by exact PSF deconvolution, via simulation with a large test set of photographic images.

1

\*\*\*\*\*

Consistent Matting for Light Field Images

Donghyeon Cho, Sunyeong Kim, and Yu-Wing Tai

Korea Advanced Institute of Science and Technology (KAIST)

Abstract. We present a new image matting algorithm to extract consistent alpha mattes across sub-images of a light field image. Instead of matting each sub-image individually, our approach utilizes the epipolar plane image (EPI) to construct comprehensive foreground and background sample sets across the sub-images without missing a true sample. The sample sets represent all color variation of foreground and background in a light field image, and the optimal alpha matte is obtained by choosing the best combination of foreground and background samples that minimizes the linear composite error subject to the EPI correspondence constraint. To further preserve consistency of the estimated alpha mattes across different sub-images, we impose a smoothness constraint along the EPI of alpha mattes. In experimental evaluations, we have created a dataset where the ground truth alpha mattes of light field images were obtained by using the blue screen technique. A variety of experiments show that our proposed algorithm produces both visually and quantitatively high-quality matting results for light field images.

Keywords: Image Matting, Light field image, EPI.

1

\*\*\*\*\*

Consensus of Regression for Occlusion-Robust

Facial Feature Localization

Xiang Yu<sup>1</sup>, Zheli Ni<sup>2</sup>, Jonathan Brandt<sup>2</sup>, and Dimitris N. Metaxas<sup>1</sup>

<sup>1</sup>Rutgers University, Piscataway, NJ 08854, USA

2Adobe Research, San Jose, CA 95110, USA

**Abstract.** We address the problem of robust facial feature localization in the presence of occlusions, which remains a lingering problem in facial analysis despite intensive long-term studies. Recently, regression-based approaches to localization have produced accurate results in many cases, yet are still subject to significant error when portions of the face are occluded. To overcome this weakness, we propose an occlusion-robust regression method by forming a consensus from estimates arising from a set of occlusion-specific regressors. That is, each regressor is trained to estimate facial feature locations under the precondition that a particular predefined region of the face is occluded. The predictions from each regressor are robustly merged using a Bayesian model that models each regressor's prediction correctness likelihood based on local appearance and consistency with other regressors with overlapping occlusion regions. After localization, the occlusion state for each landmark point is estimated using a Gaussian MRF semi-supervised learning method. Experiments on both non-occluded and occluded face databases demonstrate that our approach achieves consistently better results over state-of-the-art methods for facial landmark localization and occlusion detection.

**Keywords:** Facial feature localization, Consensus of Regression, Occlusion detection, Face alignment.

1

\*\*\*\*\*

Learning the Face Prior

for Bayesian Face Recognition

Chaochao Lu and Xiaoou Tang

Department of Information Engineering,

The Chinese University of Hong Kong, China

**Abstract.** For the traditional Bayesian face recognition methods, a simple prior on face representation cannot cover large variations in facial poses, illuminations, expressions, aging, and occlusions in the wild. In this paper, we propose a new approach to learn the face prior for Bayesian face recognition. First, we extend Manifold Relevance Determination to learn the identity subspace for each individual automatically. Based on the structure of the learned identity subspaces, we then propose to estimate Gaussian mixture densities in the observation space with Gaussian process regression. During the training of our approach, the leave-set-out algorithm is also developed for overfitting avoidance. On extensive experimental evaluations, the learned face prior can improve the performance of the traditional Bayesian face and other related methods significantly. It is also proved that the simple Bayesian face method with the learned face prior can handle the complex intra-personal variations such as large poses and large occlusions. Experiments on the challenging LFW benchmark shows that our algorithm outperforms most of the state-of-art methods.

1

\*\*\*\*\*

Spatio-temporal Event Classification Using

Time-Series Kernel Based Structured Sparsity

L'aszl'o A. Jeni<sup>1</sup>, Andr'as L'orincz<sup>2</sup>, Zolt'an Szab'o<sup>3</sup>,

Jeffrey F. Cohn<sup>1,4</sup>, and Takeo Kanade<sup>1</sup>

<sup>1</sup>Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA

<sup>2</sup>Faculty of Informatics, E'otv'os Lor'and University, Budapest, Hungary

<sup>3</sup>Gatsby Computational Neuroscience Unit, University College London, London, UK

<sup>4</sup>Department of Psychology, University of Pittsburgh, Pittsburgh, PA, USA

laszlo.jeni@ieee.org, andras.lorincz@elte.hu,

zoltan.szabo@gatsby.ucl.ac.uk, {jeffcohn, tk}@cs.cmu.edu

**Abstract.** In many behavioral domains, such as facial expression and gesture, sparse structure is prevalent. This sparsity would be well suited for event detection but for one problem. Features typically are confounded

by alignment error in space and time. As a consequence, high-dimensional representations such as SIFT and Gabor features have been favored despite their much greater computational cost and potential loss of information. We propose a Kernel Structured Sparsity (KSS) method that can handle both the temporal alignment problem and the structured sparse reconstruction within a common framework, and it can rely on simple features. We characterize spatio-temporal events as time-series of motion patterns and by utilizing time-series kernels we apply standard structured-sparse coding techniques to tackle this important problem. We evaluated the KSS method using both gesture and facial expression datasets that include spontaneous behavior and differ in degree of difficulty and type of ground truth coding. KSS outperformed both sparse and non-sparse methods that utilize complex image features and their temporal extensions. In the case of early facial event classification KSS had 10% higher accuracy as measured by F1 score over kernel SVM methods.

Keywords: structured sparsity, time-series kernels, facial expression classification, gesture recognition.

1

\*\*\*\*\*

Feature Disentangling Machine - A Novel  
Approach of Feature Selection and Disentangling  
in Facial Expression Analysis

Ping Li<sup>1</sup>, Joey Tianyi Zhou<sup>2</sup>, Ivor Wai-Hung Tsang<sup>3</sup>, Zibo Meng<sup>1</sup>,  
Shizhong Han<sup>1</sup>, and Yongdong Zhang<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of South Carolina, USA

<sup>2</sup>Center for Computational Intelligence, Nanyang Technology University, Singapore

<sup>3</sup>Center for Quantum Computation and Intelligent Systems,  
University of Technology, Australia

Abstract. Studies in psychology show that not all facial regions are of importance in recognizing facial expressions and different facial regions make different contributions in various facial expressions. Motivated by this, a novel framework, named Feature Disentangling Machine (FDM), is proposed to effectively select active features characterizing facial expressions. More importantly, the FDM aims to disentangle these selected features into non-overlapped groups, in particular, common features that are shared across different expressions and expression-specific features that are discriminative only for a target expression. Specifically, the FDM integrates sparse support vector machine and multi-task learning in a unified framework, where a novel loss function and a set of constraints are formulated to precisely control the sparsity and naturally disentangle active features. Extensive experiments on two well-known facial expression databases have demonstrated that the FDM outperforms the state-of-the-art methods for facial expression analysis. More importantly, the FDM achieves an impressive performance in a cross-database validation, which demonstrates the generalization capability of the selected features.

1

\*\*\*\*\*

Joint Unsupervised Face Alignment  
and Behaviour Analysis

Lazaros Zafeiriou, Epameinondas Antonakos,  
Stefanos Zafeiriou, and Maja Pantic

Computing Department, Imperial College London, UK

{l.zafeiriou<sup>1</sup>, e.antonakos<sup>2</sup>, s.zafeiriou<sup>1</sup>, m.pantic<sup>3</sup>}@imperial.ac.uk

Abstract. The predominant strategy for facial expressions analysis and temporal analysis of facial events is the following: a generic facial landmark tracker, usually trained on thousands of carefully annotated examples, is applied to track the landmark points, and then an analysis is performed using mostly the shape and more rarely the facial texture. This paper challenges the above framework by showing that it is feasible to perform joint

landmarks localization (i.e. spatial alignment) and temporal analysis of behaviouralsequencewiththeuseofasimplefacedetectorandasimpleshape model. To do so, we propose a new component analysis technique, which we call Autoregressive Component Analysis (ARCA), and we show how the parameters of a motion model can be jointly retrieved. The method does not require the use of any sophisticated landmark tracking methodology and simply employs pixel intensities for the texture representation. Keywords: Face alignment, time series alignment, slow feature analysis.

1

\*\*\*\*\*

## Learning a Deep Convolutional Network for Image Super-Resolution

Chao Dong<sup>1</sup>, Chen Change Loy<sup>1</sup>, Kaiming He<sup>2</sup>, and Xiaoou Tang<sup>1</sup>

<sup>1</sup>Department of Information Engineering,  
The Chinese University of Hong Kong, China

<sup>2</sup>Microsoft Research Asia, Beijing, China

Abstract. We propose a deep learning method for single image super-resolution (SR). Our method directly learns an end-to-end mapping between the low/high-resolution images. The mapping is represented as a deep convolutional neural network (CNN) [15] that takes the low-resolution image as the input and outputs the high-resolution one. We further show that traditional sparse-coding-based SR methods can also be viewed as a deep convolutional network. But unlike traditional methods that handle each component separately, our method jointly optimizes all layers. Our deep CNN has a lightweight structure, yet demonstrates state-of-the-art restoration quality, and achieves fast speed for practical on-line usage.

Keywords: Super-resolution, deep convolutional neural networks.

1

\*\*\*\*\*

## Discriminative Indexing for Probabilistic Image Patch Priors

Yan Wang<sup>1</sup>, Sunghyun Cho<sup>2</sup>, Jue Wang<sup>2</sup>, and Shih-Fu Chang<sup>1</sup>

<sup>1</sup>Dept. of Electrical Engineering, Columbia University, USA  
{yanwang,sfchang}@ee.columbia.edu

<sup>2</sup>Adobe Research, USA  
sodomau@postech.ac.kr, juewang@adobe.com

Abstract. Newly emerged probabilistic image patch priors, such as Expected Patch Log-Likelihood (EPLL), have shown excellent performance on image restoration tasks, especially deconvolution, due to its rich expressiveness. However, its applicability is limited by the heavy computation involved in the associated optimization process. Inspired by the recent advances on using regression trees to index priors defined on a Conditional Random Field, we propose a novel discriminative indexing approach on patch-based priors to expedite the optimization process. Specifically, we propose an efficient tree indexing structure for EPLL, and overcome its training tractability challenges in high-dimensional spaces by utilizing special structures of the prior. Experimental results show that our approach accelerates state-of-the-art EPLL-based deconvolution methods by up to 40 times, with very little quality compromise.

1

\*\*\*\*\*

## Modeling Video Dynamics with Deep Dynencoder

Xing Yan, Hong Chang, Shiguang Shan, and Xilin Chen

Key Lab of Intelligent Information Processing of Chinese Academy of Sciences (CAS),

Institute of Computing Technology, CAS, Beijing, 100190, China  
{xing.yan,hong.chang,shiguang.shan,xilin.chen}@vipl.ict.ac.cn

Abstract. Videos always exhibit various pattern motions, which can

be modeled according to dynamics between adjacent frames. Previous methods based on linear dynamic system can model dynamic textures but have limited capacity of representing sophisticated nonlinear dynamics. Inspired by the nonlinear expression power of deep autoencoders, we propose a novel model named dynencoder which has an autoencoder at the bottom and a variant of it at the top (named as dynpredictor). It generates hidden states from raw pixel inputs via the autoencoder and then encodes the dynamic of state transition over time via the dynpredictor. Deep dynencoder can be constructed by proper stacking strategy and trained by layer-wise pre-training and joint fine-tuning. Experiments verify that our model can describe sophisticated video dynamics and synthesize endless video texture sequences with high visual quality. We also design classification and clustering methods based on our model and demonstrate the efficacy of them on traffic scene classification and motion segmentation. ...

Keywords: Video Dynamics, Deep Model, Autoencoder, Time Series, Dynamic Textures.

1

\*\*\*\*\*

Good Image Priors for Non-blind Deconvolution:

Generic vs. Specific

Libin Sun<sup>1</sup>, Sunghyun Cho<sup>2</sup>, Jue Wang<sup>2</sup>, and James Hays<sup>1</sup>

<sup>1</sup>Brown University, Providence, RI 02912, USA

<sup>2</sup>Adobe Research, Seattle, WA 98103, USA

{lbsun,hays}@cs.brown.edu, sodomau@postech.ac.kr, juewang@adobe.com

Abstract. Most image restoration techniques build "universal" image priors, trained on a variety of scenes, which can guide the restoration of any image. But what if we have more specific training examples, e.g. sharp images of similar scenes? Surprisingly, state-of-the-art image priors don't seem to benefit from context-specific training examples. Re-training generic image priors using ideal sharp example images provides minimal improvement in non-blind deconvolution. To help understand this phenomenon we explore non-blind deblurring performance over a broad spectrum of training image scenarios. We discover two strategies that become beneficial as example images become more context-appropriate: (1) locally adapted priors trained from region level correspondences significantly outperform globally trained priors, and (2) a novel multi-scale patch-pyramid formulation is more successful at transferring mid and high frequency details from example scenes. Combining these two key strategies we can qualitatively and quantitatively outperform leading generic non-blind deconvolution methods when context-appropriate example images are available. We also compare to recent work which, like ours, tries to make use of context-specific examples.

Keywords: deblur, non-blind deconvolution, gaussian mixtures, image pyramid, image priors, camera shake.

1

\*\*\*\*\*

Image Deconvolution Ringing Artifact Detection

and Removal via PSF Frequency Analysis

Ali Mosleh<sup>1</sup>, J.M. Pierre Langlois<sup>1</sup>, and Paul Green<sup>2</sup>

<sup>1</sup>Ecole Polytechnique de Montr  al, Canada

<sup>2</sup>Algolux, Canada

{ali.mosleh,pierre.langlois}@polymtl.ca, paul.green@algolux.com

Abstract. We present a new method to detect and remove ringing artifacts produced by the deconvolution process in image deblurring techniques. The method takes into account non-invertible frequency components of the blur kernel used in the deconvolution. Efficient Gabor wavelets are produced for each non-invertible frequency and applied on the deblurred image to generate a set of filter responses that reveal existing ringing artifacts. The set of Gabor filters is then employed in a



regularization scheme to remove the corresponding artifacts from the deblurred image. The regularization scheme minimizes the responses of the reconstructed image to these Gabor filters through an alternating algorithm in order to suppress the artifacts. As a result of these steps we are able to significantly enhance the quality of the deblurred images produced by deconvolution algorithms. Our numerical evaluations using a ringing artifact metric indicate the effectiveness of the proposed deringing method.

Keywords: deconvolution, image deblurring, point spread function, ringing artifacts, zero-magnitude frequency.

1

\*\*\*\*\*

#### View-Consistent 3D Scene Flow Estimation over Multiple Frames

Christoph Vogell, Stefan Roth, and Konrad Schindler<sup>1</sup>

<sup>1</sup>Photogrammetry and Remote Sensing, ETH Zurich, Switzerland

<sup>2</sup>Department of Computer Science, TU Darmstadt, Germany

Abstract. We propose a method to recover dense 3D scene flow from stereo video. The method estimates the depth and 3D motion field of a dynamic scene from multiple consecutive frames in a sliding temporal window, such that the estimate is consistent across both viewpoints of all frames within the window. The observed scene is modeled as a collection of planar patches that are consistent across views, each undergoing a rigid motion that is approximately constant over time. Finding the patches and their motions is cast as minimization of an energy function over the continuous plane and motion parameters and the discrete pixel-to-plane assignment. We show that such a view-consistent multi-frame scheme greatly improves scene flow computation in the presence of occlusions, and increases its robustness against adverse imaging conditions, such as specularities. Our method currently achieves leading performance on the KITTI benchmark, for both flow and stereo.

1

\*\*\*\*\*

#### Hand Waving Away Scale

Christopher Ham, Simon Lucy<sup>2</sup>, and Surya Singh<sup>1</sup>

<sup>1</sup>Robotics Design Lab, The University of Queensland, Australia

<sup>2</sup>Robotics Institute, Carnegie Mellon University, USA

{c.ham, spns}@uq.edu.au, slucey@cs.cmu.edu

Abstract. This paper presents a novel solution to the metric reconstruction of objects using any smart device equipped with a camera and an inertial measurement unit (IMU). We propose a batch, vision centric approach which only uses the IMU to estimate the metric scale of a scene reconstructed by any algorithm with Structure from Motion like (SfM) output. IMUs have a rich history of being combined with monocular vision for robotic navigation and odometry applications. These IMUs require sophisticated and quite expensive hardware rigs to perform well. IMUs in smart devices, however, are chosen for enhancing interactivity - a task which is more forgiving to noise in the measurements. We anticipate, however, that the ubiquity of these "noisy" IMUs makes them increasingly useful in modern computer vision algorithms. Indeed, we show in this work how an IMU from a smart device can help a face tracker to measure pupil distance, and an SfM algorithm to measure the metric size of objects. We also identify motions that produce better results, and develop a heuristic for estimating, in real-time,

when enough data has been collected for an accurate scale estimation.

Keywords: Smart devices, IMU, metric, 3D reconstruction.

1

\*\*\*\*\*

#### A Non-Linear Filter for Gyroscope-Based Video Stabilization

Steven Bell<sup>1</sup>, Alejandro Troccoli<sup>2</sup>, and Kari Pulli<sup>2</sup>

1Stanford University, Stanford, CA, USA  
sebell@stanford.edu

2NVIDIA Research, Santa Clara, CA, USA  
{atroccoli,karip}@nvidia.com

Abstract. We present a method for video stabilization and rolling-shutter correction for videos captured on mobile devices. The method uses the data from an on-board gyroscope to track the camera's angular velocity, and can run in real time within the camera capture pipeline. We remove small motions and rolling-shutter distortions due to hand shake, creating the impression of a video shot on a tripod. For larger motions, we alter the camera's angular velocity to produce a smooth output. To meet the latency constraints of a real-time camera capture pipeline, our alter operates on a small temporal window of three to five frames. Our algorithm performs better than the previous work that uses a gyroscope to stabilize a video stream, and at a similar level with respect to current feature-based methods.

Keywords: video stabilization, rolling-shutter, gyroscopes.

1

\*\*\*\*\*

Multi-modal and Multi-spectral Registration  
for Natural Images

Xiaoyong Shen<sup>1</sup>, Li Xu<sup>2</sup>, Qi Zhang<sup>1</sup>, and Jia-Ya Jiang<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong, China

<sup>2</sup>Image & Visual Computing Lab, Lenovo R&T,

Project Website, Hong Kong, China

<http://www.cse.cuhk.edu.hk/leojia/projects/multimodal>

Abstract. Images now come in different forms - color, near-infrared, depth, etc. - due to the development of special and powerful cameras in computer vision and computational photography. Their cross-modal correspondence establishment is however left behind. We address this challenging dense matching problem considering structure variation possibly existing in these image sets and introduce new model and solution. Our main contribution includes designing the descriptor named robust selective normalized cross correlation (RSNCC) to establish dense pixel correspondence in input images and proposing its mathematical parameterization to make optimization tractable. A computationally robust framework including global and local matching phases is also established. We build a multi-modal dataset including natural images with labeled sparse correspondence. Our method will benefit image and vision applications that require accurate image alignment.

Keywords: multi-modal, multi-spectral, dense matching, variational model.

1

\*\*\*\*\*

Using Isometry to Classify

Correct/Incorrect 3D -2D Correspondences

Toby Collins and Adrien Bartoli

ALCoV-ISIT, UMR 6284 CNRS/Universit e d'Auvergne, Clermont-Ferrand, France

Abstract. Template-based methods have been successfully used for surface detection and 3D reconstruction from a 2D input image, especially when the surface is known to deform isometrically. However, almost all such methods require that keypoint correspondences be first matched between the template and the input image. Matching thus exists as a current limitation because existing methods are either slow or tend to perform poorly for discontinuous or unsmooth surfaces or deformations. This is partly because the 3D isometric deformation constraint cannot be easily used in the 2D image directly. We propose to resolve that difficulty by detecting incorrect correspondences using the isometry constraint directly in 3D. We do this by embedding a set of putative correspondences in 3D space, by estimating their depth and local 3D orientation in the

input image, from local image warps computed quickly and accurately by means of Inverse Composition. We then relax isometry to inextensibility to get a first correct/incorrect classification using simple pairwise constraints. This classification is then efficiently refined using higher-order constraints, which we formulate as the consistency between the correspondences' local 3D geometry. Our algorithm is fast and has only one free parameter governing the precision/recall trade-off. We show experimentally that it significantly outperforms state-of-the-art.

1

\*\*\*\*\*

Bilateral Functions for Global Motion Modeling

Wen-Yan Daniel Lin<sup>1</sup>, Ming-Ming Cheng<sup>2</sup>, Jiangbo Lu<sup>1</sup>, Hongsheng Yang<sup>3</sup>,  
Minh N. Do<sup>4</sup>, and Philip Torr<sup>2</sup>

<sup>1</sup>Advanced Digital Sciences Center, Singapore

<sup>2</sup>Oxford University, UK

<sup>3</sup>University of North Carolina at Chapel Hill, USA

<sup>4</sup>University of Illinois at Urbana-Champaign, USA

Abstract. This paper proposes modeling motion in a bilateral domain that augments spatial information with the motion itself. We use the bilateral domain to reformulate a piecewise smooth constraint as continuous global modeling constraint. The resultant model can be robustly computed from highly noisy scattered

feature points using a global minimization. We demonstrate how the model can reliably obtain large numbers of good quality correspondences over wide baselines, while keeping outliers to a minimum.

1

\*\*\*\*\*

VCDB: A Large-Scale Database

for Partial Copy Detection in Videos

Yu-Gang Jiang, Yudong Jiang, and Jiajun Wang

School of Computer Science, Shanghai Key Laboratory of Intelligent  
Information Processing, Fudan University, Shanghai, China

ygj@fudan.edu.cn

Abstract. The task of partial copy detection in videos aims at finding if one or more segments of a query video have (transformed) copies in a large dataset. Since collecting and annotating large datasets of real partial copies are extremely time-consuming, previous video copy detection research used either small-scale datasets or large datasets with simulated partial copies by imposing several pre-defined transformations (e.g., photometric or geometric changes). While the simulated datasets were useful for research, it is unknown how well the techniques developed on such data work on real copies, which are often too complex to be simulated. In this paper, we introduce a large-scale video copy database (VCDB) with over 100,000 Web videos, containing more than 9,000 copied segment pairs found through careful manual annotation. We further benchmark a baseline system on VCDB, which has demonstrated state-of-the-art results in recent copy detection research. Our evaluation suggests that existing techniques—which have shown near-perfect results on the simulated benchmarks—are far from satisfactory in detecting complex real copies. We believe that the release of VCDB will largely advance the research around this challenging problem.

Keywords: Video copy detection, benchmark dataset, frame matching, temporal alignment.

1

\*\*\*\*\*

Single-Image Super-Resolution: A Benchmark

Chih-Yuan Yang<sup>1</sup>, Chao Ma<sup>1,2</sup>, and Ming-Hsuan Yang<sup>1</sup>

<sup>1</sup>University of California at Merced, USA

<sup>2</sup>Shanghai Jiao Tong University, China

{cyang35, cma26, mhyang}@ucmerced.edu

Abstract. Single-image super-resolution is of great importance for vision applications, and numerous algorithms have been proposed in recent years. Despite the demonstrated success, these results are often generated based on different assumptions using different datasets and metrics. In this paper, we present a systematic benchmark evaluation for state-of-the-art single-image super-resolution algorithms. In addition to quantitative evaluations based on conventional full-reference metrics, human subject studies are carried out to evaluate image quality based on visual perception. The benchmark evaluations demonstrate the performance and limitations of state-of-the-art algorithms which sheds light on future research in single-image super-resolution.

Keywords: Single-image super-resolution, performance evaluation, metrics, Gaussian blur kernel width.

1

\*\*\*\*\*

Well Begun Is Half Done:

Generating High-Quality Seeds

for Automatic Image Dataset Construction from Web

Yan Xia<sup>1</sup>, Xudong Cao<sup>2</sup>, Fang Wen<sup>2</sup>, and Jin Sun<sup>2</sup>

<sup>1</sup>University of Science and Technology of China

<sup>2</sup>Microsoft Research Asia, Beijing, China

Abstract. We present a fully automatic approach to construct a large-scale, high-

precision dataset from noisy web images. Within the entire pipeline, we focus on generating high quality seed images for subsequent dataset growing. High quality seeds are essential as we revealed, but they have received relatively less attention

in previous works with respect to how to automatically generate them. In this work, we propose a density score based on rank-order distance to identify positive seed images. The basic idea is images relevant to a concept typically are tightly

clustered, while the outliers are widely scattered. Through adaptive thresholding,

we guarantee the selected seeds as numerous and accurate as possible. Starting with the high quality seeds, we grow a high quality dataset by dividing seeds and conducting iterative negative and positive mining. Our system can automatically collect thousands of images for one concept/class, with a precision rate of 95% or more. Comparisons with recent state-of-the-arts also demonstrate our method's superior performance.

1

\*\*\*\*\*

Zero-Shot Learning via Visual Abstraction

Stanislaw Antol<sup>1</sup>, C. Lawrence Zitnick<sup>2</sup>, and Devi Parikh<sup>1</sup>

<sup>1</sup>Virginia Tech, Blacksburg, VA, USA

<sup>2</sup>Microsoft Research, Redmond, WA, USA

Abstract. One of the main challenges in learning fine-grained visual categories is gathering training images. Recent work in Zero-Shot Learning (ZSL) circumvents this challenge by describing categories via attributes or text. However, not all visual concepts, e.g., two people dancing, are easily amenable to such descriptions. In this paper, we propose a new modality for ZSL using visual abstraction to learn difficult-to-describe concepts. Specifically, we explore concepts related to people and their interactions with others. Our proposed modality allows one to provide training data by manipulating abstract visualizations, e.g., one can illustrate interactions between two clipart people by manipulating each person's pose, expression, gaze, and gender. The feasibility of our approach is shown on a human pose dataset and a new dataset containing complex interactions between two people, where we outperform several baselines. To better match across the two domains, we learn an explicit mapping between the abstract and real worlds.

Keywords: zero-shot learning, visual abstraction, synthetic data, pose.

1

\*\*\*\*\*

Zero-Shot Learning via Visual Abstraction 403

Moreover, our easy-to-use interface results in biases in the illustrations ( e.g .,

the interface does not allow for out-of-plane rotation). To account for these human tendencies, as well as interface biases, we learn an explicit mapping from the features extracted from illustrations to the features extracted from real images.

This allows us to improve performance on instance-level ZSL. Our visual abstraction interface, code, and datasets are publicly available.

## 2 Related Work

We discuss existing work on zero-shot learning, learning with synthetic data, learning semantic relations, pose estimation, and action recognition.

Zero-Shot Learning (ZSL): The problem of learning models of visual concepts without example images of the concepts is called Zero-Shot Learning. Attributes (mid-level, visual, and semantic features) [9, 10, 15, 16] provide a natural interface for ZSL [16], where an unseen class is described by a list of attributes. Equipped

with a set of pre-trained attribute classifiers, a test image can be probabilistically matched to each of these attribute signatures and be classified as the category with the highest probability. Instead of using a list of attributes, recent work [7] has leveraged more general textual descriptions of categories to build visual

models of these categories. Our work takes a fundamentally different approach to ZSL. We propose a strictly visual modality to allow a supervisor to train a model for visual concepts that may not be easily describable in semantic terms, e.g., poses of people, interactions between people.

Learning With Synthetic Data: Our work introduces the use of abstract visualizations as a modality to train visual models in a ZSL setting. Previously,

papers have explored the use of synthetic data to aid in the training of vision algorithms.

In many object recognition tasks, it is common to perturb the training data using affine warps to augment the training data [14]. Computer-generated scenes may also be used to evaluate recognition systems [13]. Shotton et al. [23]

used synthetically generated depth data depicting humans to learn a human pose detector from this depth data. Unlike these approaches, we are trying to learn high-level, complex concepts where it is not feasible to automatically generate synthetic data, so we must rely on humans to create our synthetic data. Most similar to our work, the problem of semantic scene understanding using abstract scenes was studied in [31]. They use a dataset of simple sentences corresponding

to abstract scenes to learn a mapping from sentences to abstract scenes. Recently, sequences of abstract scenes were used to predict which objects will move in the near future [11]. Unlike these works, we use abstraction to learn visual models that can be applied to real images. Sketch-based image retrieval [6, 24] allows users to search for an image by sketching the concept. Sketching complex interactions between people would be time consuming, and likely inaccurate for most lay users. More importantly, our modality has the potential to augment the abstract scenes with a large variety of visual cues ( e.g., gender, ethnicity, clothing, background) that would be cumbersome for users to convey via sketches

\*\*\*\*\*

Zero-Shot Learning via Visual Abstraction 405

(e.g., somebody's part can be cropped out), which makes this dataset challenging (e.g., right side of Figure 1). This resulted in 3,600 initial images.

Real Image Annotations: We also used AMT to collect various image annotations.

notations that are needed for our features via different custom interfaces. The pose annotation interface prompted the worker with one of our images and its corresponding sentence. We highlighted whether the worker should be annotating Person A or Person B in the sentence. The worker annotates the person's 14 body parts (right side of Figure 3). The worker provides their best guess if the part is occluded and responds "not present" if it is not within the image border.

We had 5 workers annotate each person in each image and averaged them for the final ground truth pose annotations. In addition, workers annotated ground truth eyegaze (i.e., looking to the image left or right), facial expression (i.e., one of six prototypic emotional expressions [5] plus a neutral expression), and gender of each person via separate interfaces. We selected the mode of their responses for our final annotation. In addition to collecting the annotation of interest, two interfaces asked one additional question each. One asked if the prompted image contained exactly two main people or not and the other asked if the annotated pose overlaid on the prompted image was of good quality or not. We used the last two annotation queries to remove poor quality work. Additionally, a GIST-based [20] image matching scheme was used to remove duplicates. Removing these images gave us our final annotated dataset with 3,172 images (52.9 images per category on average). Some examples can be found in the bottom part of Figure 1 and the rightmost two columns of Figure 5. More details about our interfaces and our procedure can be found in the supplementary material.

### 3.2 PARSE

We also use a subset of the standard PARSE [21] dataset, which originally contains 305 images of individuals in various poses. We created a list of categories

that frequently appear in the PARSE dataset (e.g., "is dunking," "is diving for an object"). From the images that belong to these categories, we removed those that were used to train the pose detector [26]. Some categories (e.g., "is standing-

ing") had disproportionately large number of images, so we removed images at random from these categories. This leaves us with 108 images in our dataset (7.7 images per category on average). We also collected the same annotations as in Section 3.1, except for pose (since ground truth pose annotations are already available with the dataset). See the supplementary material for more details.

## 4 Our Approach

In this section, we present our new modality for ZSL. We begin by introducing our user interface for collecting visual illustrations for training. We then describe

the novel features that are extracted from our abstract illustrations and real images. Finally, we describe the approach used to train our models. The results of various experiments follow in Section 5

\*\*\*\*\*

406 S. Antol, C.L. Zitnick, and D. Parikh

Fig. 2. User interface (with random initialization) used to collect abstract illustrations

on AMT. Workers were able to manipulate pose, expression, gaze direction, and gender.

### 4.1 Visual Abstraction Interface

For our domain of interest, we conjectured that our concepts depend primarily on four main factors: pose, eye gaze, facial expression, and gender. Some other factors that we do not model, but may also be important are clothing, the presence of other objects, and scene context. A screenshot of our user interface is shown in Figure 2. Initially, two people (one blond-haired and one brown-haired) are shown with random poses, gaze directions (i.e., "lip"), expressions

, and genders. We allow our subjects to continuously manipulate the poses (i.e., joint angles and positions) of both people by dragging on the various body parts. They may horizontally lip the people to change their perceived eye gaze

direction. The facial expressions are chosen from the same selection as is used for the annotation of real images (Section 3.1). Finally, the subjects may select one of the two predominant genders for each clipart person.

To collect our training data for category-level ZSL, we prompt the user with a sentence to illustrate using the interface (e.g., "Person A is dancing with Person B.", "A person is dunking."). To promote diversity, we encouraged them to imagine any objects or background, as long as the poses are consistent with the imagined scene (e.g., a worker can imagine a chair and illustrate someone sitting on it). The interface includes buttons to annotate which clipart person corresponds to which person in the sentence. Some illustrations are shown on the left side of Figure 1 and in the left three columns of Figure 5. For the PARSE concepts, the interface is the same except that only one person is present. For instance-level ZSL, we modify our previous interface. Instead of sentences, we first (briefly, for 2 seconds) show the user a real image and then they recreate

it (from memory) as best they can. The stated goal is to recreate the real image so another person would be able to select the shown image from a collection of real images. This mimics the scenario when a person is searching for a specific image: they might be clear on the semantically important aspects while having a fuzzier or skewed notion of other aspects. Another bias of the illustrations occurs

when it is impossible to recreate the real image exactly due to the limitations of

the interface, such as not being able to change the height of the clipart people, the interface not allowing for out-of-plane rotation, etc

\*\*\*\*\*

Zero-Shot Learning via Visual Abstraction 407

## 4.2 Relation and Appearance Features

Using the annotations described in Section 3.1 (i.e., pose, gaze, expression, and

gender) for persons denoted by  $i$  and  $j$ , we compute a set of relation and appearance features. Some of our relation features are distance-based and some are angle-based. All distance-based features use Gaussians placed at different positions to capture relative distance. The Gaussians'  $\sigma$  parameters are proportional to the scale of each person. A person's scale is defined as the distance between their head and the center of their shoulders and hips. Unless otherwise noted, all angles/orientations are w.r.t. the image frame's x-axis. They are represented

by 12 dimensional unit histograms with each bin corresponding to  $\pi/6$  radians. Soft assignments are made to the histograms using linear weighting. The first two sets of features, Basic and Gaze, account for both people. The remaining five feature sets are described for a single person and must be evaluated twice (swapping  $i$  and  $j$ ) and concatenated. The feature sets are described below.

**Basic:** This feature set encodes basic relation properties between two people, such as relative orientation and distance. We calculate each person's body angle (in the image frame). This is calculated from the image coordinates for the head and mid-point between shoulders. We place Gaussians at the center of the people and then use the distance between them to evaluate the Gaussian functions. We also calculate the angle (in the image frame) between the centers of the two people. This gives us a total of  $2 * (12 + 1) + 12 = 38$  features. They can be thought of as simplifying the people into two boxes (possibly having different scale parameters) with certain orientations and looking at the relative position and angle between their centers.

**Gaze:** The gaze feature set is encoded using 5 binary features, corresponding to looking at  $j$ ,  $j$  looking at  $i$ , both people are looking at each other, both people are looking away from each other, and both people are looking in the same direction. To determine if  $i$  is looking at  $j$ , we check if  $j$ 's neck is in the

appropriate region of the image. The image is divided into two parts by extending the line between  $i$ 's head and neck and the appropriate region is defined to be

the area where  $i$  is looking (which depends on  $i$ 's gaze direction). Once we have both looking at and vice versa features, we compute the remaining three gaze features via the appropriate logic operations (e.g., if  $i$  is looking at  $j$  and  $j$  is looking at  $i$ , then the looking-at-each-other feature is true).

**Global:** This feature set encodes the general position of the joints in reference to a body. Three Gaussians are placed in a  $3 \times 1$  grid on the image based on the body's size and orientation (the blue circles in Figure 3). The positions of one person's 8 joints (two for each limb) are evaluated using all Gaussians from both Gaussian sets (

i.e., person  $i$ 's joints relative to person  $i$ 's global Gaussians and person  $j$ 's global Gaussians), giving us a total of  $8 \times 3 \times 2 = 48$  features.

**Contact:** This feature set encodes the specific location of the joints in reference

to other body parts. For each person, we place Gaussians at 13 positions:

\*\*\*\*\*

Zero-Shot Learning via Visual Abstraction 409

if an image represents a specific concept, i.e., given a test real image, we wish to determine which specific abstract visualization (instance) corresponds to the real image. For this, we use Nearest Neighbor matching. Since our features are from two different domains, learning a mapping between them could improve the matching performance. This is described next.

#### 4.4 Mapping From Abstract to Real for the Instance-level Model

We learn a mapping between the domain of abstract images and the domain of real images. To learn such a mapping, we need examples that correspond to the same thing in both domains. We use some of our instance-level illustrations (Section 4.1) as these abstract-real pairs. The mapping can learn to correct for

both user and interface biases discussed in Section 4.1.

Simpler techniques, such as Canonical Correspondence Analysis [12], did not learn a good mapping between the abstract and real worlds. We found that General Regression Neural Networks (GRNN) [25] did better. We also found that converting from our abstract features into "real" features performed better than converting real features into "abstract" features. Thus, the GRNN's input is all of the abstract features and its output is all of the real features.

### 5 Experimental Results

In this section, we describe our experiments which show that our new modality for ZSL is able to create models that can learn category-level (Section 5.1) and

instance-level (Section 5.2) visual concepts. We perform an ablation study on different feature sets, showing their performance contribution (Section 5.3). Finally, we utilize a state-of-the-art pose detector on both INTERACT and PARSE datasets to investigate our approach in a more automatic setting (Section 5.4).

#### 5.1 Category-Level Zero-Shot Learning

We begin by experimenting with the ability of our novel modality to learn our category-level concepts, i.e., classifying images into one of the semantic descriptions,

such as "A is kicking B." To acquire the required training illustrations, we ran our visual abstraction interface with sentence prompts (described in Section 4.1) on AMT. We had 50 workers create an abstract illustration for each of the 60 semantic concepts from INTERACT (Section 3.1) and the 14 semantic concepts from PARSE (Section 3.2). After removing poor quality work, we are left with 3,000 and 696 illustrations, respectively.

The setup for all category-level ZSL experiments (unless otherwise noted) is described here. Using the abstract illustrations, we train multiple one-vs-all linear SVMs (liblinear [8]) with the cost parameter,  $C$ , set to 0.01, which worked

reasonably well across all experiments. For INTERACT, there is ambiguity (at test time) as to which person is Person A and which person is Person B. To account for this, we evaluate each of the classifiers using both orderings, select the

\*\*\*\*\*



## 5.2 Instance-Level Zero-Shot Learning

We also test the ability of our new modality to learn instance-level concepts. To acquire the necessary training illustrations, we ran our visual abstraction interface with image prompts (as described in Section 4.1) on AMT. We showed a real image (one of 3,172 from INTERACT and one of 305 from PARSE) for two seconds to the workers, who recreated it using the interface. Through a pilot

study, just as in [6], we found two seconds to be sufficient for people to capture the more salient aspects of the image. It is unlikely that a user would have every detail of the instance in mind when trying to train a model for a specific concept

and we wanted to mimic this in our set up. We had 3 workers recreate each of the images, and after manually removing work from problematic workers, we are left with 8,916 and 914 illustrations for INTERACT and PARSE, respectively.

We perform classification via nearest-neighbor matching. If the real image's features match the features of any of the (up to) 3 illustration instances that workers created for it, we have found a correct label. We vary  $K$ , the number of nearest neighbors that are considered, and evaluate the percentage of real images that have a correct label within those  $K$  neighbors. We normalized  $K$  by the total number of illustrations. We need a training dataset to learn a mapping between the abstract and real worlds, i.e., training the GRNN from Section 4.4. For INTERACT, we split the categories into 39 seen categories for training and 21 unseen categories for testing to minimize learning biases specific to specific categories (i.e., verb phrases).

The results are averaged over 10 random seen/unseen category splits. For PARSE, the training data corresponds to the 197 images in Fig. 5. The left columns show 5 random illustrations (of 50) used for classifier training.

Columns 6 and 7 contain the most confident true positive and false positive for a given category, respectively. Mistakes include choosing a semantically reasonable verb (top),

choosing the incorrect preposition (middle), and incorrect prediction due to the pose

similarity between two classes (bottom). More examples are in the supplement

\*\*\*\*\*

0481216

Random

B

C

G

O

B+C

B+G

B+O

B+E+Z+S

B+C+G+O

B+C+G+O+E

C+G+O+E+Z+S

B+G+O+E+Z+S

B+C+O+E+Z+S

B+C+G+E+Z+S

B+C+G+O+Z+S

B+C+G+O+E+S

B+C+G+O+E+Z

B+C+G+O+E+Z+SPP

YR-BB

Random

1

1Mean of Class-wise  
 Raw Accuracies (%) INTERACT  
 Features:  
 (B)asic  
 (C)ontact  
 (G)lobal  
 (O)rientation  
 (E)xpression  
 Ga(Z)e  
 S for Gender

Fig. 7. We plot classification performance for INTERACT using different subsets of features. Some features, like Global, are more informative than others. Of the appearance-

based features, Expression turns out to be most informative, presumably when body pose features are similar (e.g., "wrestling" vs. "hugging").

#### 5.4 Automatic Pose Evaluation

In this section, we do an evaluation of our category-level ZSL task using the current state-of-the-art pose detector developed by Yang and Ramanan [27]. We utilized the pre-trained PARSE model and detected the pose on both the INTERACT and the PARSE datasets. For the expression, gaze, and gender features, we continue to use human annotations. These results (YR) are shown in Figures 4, 6, and 7. As expected, due to the pose detector being developed for PARSE, automatic detection on the PARSE dataset yields reasonable performance (compared to perfect pose). The results on INTERACT do not perform nearly as well, although it still outperforms the baselines. To boost the perfor-

mance of the pose detector on INTERACT, we also experimented with providing ground truth bounding boxes (YR-BB), which results in better performance. INTERACT is significantly more challenging than PARSE for automatic pose detection. Thus, it is not surprising that incorrectly detected poses confuse our

models. Properties that make INTERACT particularly challenging include: images from arbitrary perspectives, more difficult (for the detector) poses (e.g., "crawling," "lying"), overlapping people (e.g., "hugging," "standing in front of"), and incomplete poses (i.e., not all body parts are present). We investigated

this latter point by selecting images from INTERACT based on the number of parts present in the image. There are 14 parts per person and we ensure that both people have at least a certain number of parts. Requiring all parts to be within the image reduces INTERACT to 1,689 images (from 3,172). 91.5% of our images contain at least 7 parts per person. More of these details can be found in the supplementary material. We re-evaluate our category-level ZSL performance (at 50 training illustrations per category) as we vary the part threshold and show our results in Figure 8. Although there is some noise, both the perfect pose and automatic pose detection methods show an increase in accuracy as we require

\*\*\*\*\*

#### Zero-Shot Learning via Visual Abstraction 415

##### References

1. Ali, S., Shah, M.: Human action recognition in videos using kinematic features and multiple instance learning. PAMI (2010)
2. Bourdev, L., Malik, J.: Poselets: Body part detectors trained using 3d human pose annotations. In: ICCV (2009)
3. Chakraborty, I., Cheng, H., Javed, O.: 3d visual proxemics: Recognizing human interactions in 3d from a single image. In: CVPR (2013)
4. Choi, W., Shahid, K., Savarese, S.: Learning context for collective activity recognition. In: CVPR (2011)
5. Darwin, C.: The Expression of the Emotions in Man and Animals. Oxford

University Press (1998)

6. Eitz, M., Hildebrand, K., Boubekeur, T., Alexa, M.: Sketch-based image retrieval:

Benchmark and bag-of-features descriptors. *IEEE Transactions on Visualization and Computer Graphics* (2011)

7. Elhoseiny, M., Saleh, B., Elgammal, A.: Write a classifier: Zero-shot learning using

purely textual descriptions. In: *ICCV* (2013)

8. Fan, R.E., Chang, K.W., Hsieh, C.J., Wang, X.R., Lin, C.J.: LIBLINEAR:

A library for large linear classification. *JMLR* (2008)

9. Farhadi, A., Endres, I., Hoiem, D., Forsyth, D.: Describing objects by their attributes. In: *ICCV* (2009)

10. Ferrari, V., Zisserman, A.: Learning visual attributes. In: *NIPS* (2007)

11. Fouhey, D.F., Zitnick, C.L.: Predicting object dynamics in scenes. In: *CVPR* (2014)

12. Hotelling, H.: Relations between two sets of variates. *Biometrika* (1936)

13. Kaneva, B., Torralba, A., Freeman, W.T.: Evaluation of image features using a

photorealistic virtual world. In: *ICCV* (2011)

14. Krizhevsky, A., Sutskever, I., Hinton, G.: ImageNet classification with deep convolutional

neural networks. In: *NIPS* (2012)

15. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Attribute and simile classifiers

for face verification. In: *ICCV* (2009)

16. Lampert, C.H., Nickisch, H., Harmeling, S.: Learning to detect unseen object classes by

between-class attribute transfer. In: *CVPR* (2009)

17. Lan, T., Wang, Y., Yang, W., Mori, G.: Beyond actions: Discriminative models for

contextual group activities. In: *NIPS* (2010)

18. Larochelle, H., Erhan, D., Bengio, Y.: Zero-data learning of new tasks. In: *AAAI* (2008)

19. Marin-Jimenez, M., Zisserman, A., Eichner, M., Ferrari, V.: Detecting people looking

at each other in videos. *IJCV* (2013)

20. Oliva, A., Torralba, A.: Modeling the shape of the scene: A holistic representation

of the spatial envelope. *IJCV* (2001)

21. Ramanan, D.: Learning to parse images of articulated bodies. In: *NIPS* (2007)

22. Sadeghi, M.A., Farhadi, A.: Recognition using visual phrases. In: *CVPR* (2011)

23. Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore, R., Kipman, A., Blake, A.: Real-time human pose recognition in parts from single depth

images. In: *CVPR* (2011)

24. Smeulders, A.W.M., Worring, M., Santini, S., Gupta, A., Jain, R.: Content-based

image retrieval at the end of the early years. *PAMI* (2000)

25. Specht, D.F.: The general regression neural network-rediscovered. *Neural Networks*

(1993)

\*\*\*\*\*  
Discovering Groups of People in Images

Wongun Choi<sup>1</sup>, Yu-Wei Chao<sup>2</sup>, Caroline Pantofaru<sup>3</sup>, and Silvio Savarese<sup>4</sup>

<sup>1</sup>NEC Laboratories, USA

<sup>2</sup>University of Michigan, Ann Arbor, USA

<sup>3</sup>Google, Inc, USA

<sup>4</sup>Stanford University, USA

Abstract. Understanding group activities from images is an important yet challenging task. This is because there is an exponentially large number of semantic and geometrical relationships among individuals that one must model in order to effectively recognize and localize the group activities. Rather than focusing on

n di-

rectly recognizing group activities as most of the previous works do, we advocate

the importance of introducing an intermediate representation for modeling groups of humans which we call structure groups. Such groups define the way people spatially interact with each other. People might be facing each other to talk, while others sit on a bench side by side, and some might stand alone. In this paper we contribute a method for identifying and localizing these structured groups in a

single image despite their varying viewpoints, number of participants, and occlusions. We propose to learn an ensemble of discriminative interaction patterns to encode the relationships between people in 3D and introduce a novel efficient iterative augmentation algorithm for solving this complex inference problem. A nice byproduct of the inference scheme is an approximate 3D layout estimate of the structured groups in the scene. Finally, we contribute an extremely challenging new dataset that contains images each showing multiple people performing multiple activities. Extensive evaluation confirms our theoretical findings.

Keywords: Group discovery, Social interaction, Activity recognition.

1

\*\*\*\*\*

Untangling Object-View Manifold

for Multiview Recognition and Pose Estimation

Amr Bakry and Ahmed Elgammal

Department of Computer Science, Rutgers University

Piscataway, NJ, USA

Abstract. The problem of multi-view/view-invariant recognition remains one of the most fundamental challenges to the progress of the computer vision. In this paper we consider the problem of modeling the combined object-viewpoint manifold. The shape and appearance of an object in a given image is a function of its category, style within category, viewpoint, and several other factors. The visual manifold (in any chosen feature representation space) given all these variability collectively is very hard and even impossible to model. We propose an efficient computational framework that can untangle such a complex manifold, and achieve a model that separates a view-invariant category representation, from category-invariant pose representation. We outperform the state of the art in the three widely used multiview dataset, for both category recognition, and pose estimation.

1

\*\*\*\*\*

Parameterizing Object Detectors

in the Continuous Pose Space

Kun He<sup>1</sup>, Leonid Sigal<sup>2</sup>, and Stan Sclaroff<sup>1</sup>

<sup>1</sup>Computer Science Department, Boston University, USA

<sup>2</sup>Disney Research Pittsburgh, USA

{hekun,sclaroff}@cs.bu.edu, lsigal@disneyresearch.com

Abstract. Object detection and pose estimation are interdependent problems in computer vision. Many past works decouple these problems, either by discretizing the continuous pose and training pose-specific object detectors, or by building pose estimators on top of detector outputs. In this paper, we propose a structured kernel machine approach to treat object detection and pose estimation jointly in a mutually beneficial way. In our formulation, a unified, continuously parameterized, discriminative appearance model is learned over the entire pose space. We propose a cascaded discrete-continuous algorithm for efficient inference, and give effective online constraint generation strategies for learning our model using structural SVMs. On three standard benchmarks, our method performs better than, or on par with, state-of-the-art methods in the combined task of object detection and pose estimation.

Keywords: object detection, continuous pose estimation.

1

\*\*\*\*\*

## Jointly Optimizing 3D Model Fitting and Fine-Grained Classification

Yen-Liang Lin<sup>1</sup>, Vlad I. Morariu<sup>2</sup>, Winston Hsu<sup>1</sup>, and Larry S. Davis<sup>2</sup>

<sup>1</sup>National Taiwan University, Taipei, Taiwan

<sup>2</sup>University of Maryland, College Park, MD, USA

yenliang@cmlab.csie.ntu.edu.tw, whsu@ntu.edu.tw,

{morariu,lsd}@umiacs.umd.edu

**Abstract.** 3D object modeling and fine-grained classification are often treated as separate tasks. We propose to optimize 3D model fitting and fine-grained classification jointly. Detailed 3D object representations encode more information (e.g., precise part locations and viewpoint) than traditional 2D-based approaches, and can therefore improve fine-grained classification performance. Meanwhile, the predicted class label can also improve 3D model fitting accuracy, e.g., by providing more detailed class-specific shape models. We evaluate our method on a new fine-grained 3D car dataset (FG3D Car), demonstrating our method outperforms several state-of-the-art approaches. Furthermore, we also conduct a series of analyses to explore the dependence between fine-grained classification performance and 3D models.

1

\*\*\*\*\*

## Pipelining Localized Semantic Features for Fine-Grained Action Recognition

Yang Zhou<sup>1</sup>, Bingbing Ni<sup>2</sup>, Shuicheng Yan<sup>3</sup>, Pierre Moulin<sup>4</sup>, and Qi Tian<sup>1</sup>

<sup>1</sup>University of Texas at San Antonio, USA

<sup>2</sup>Advanced Digital Sciences Center, Singapore

<sup>3</sup>National University of Singapore, Singapore

<sup>4</sup>University of Illinois at Urbana-Champaign, USA

myh511@my.utsa.edu, bingbing.ni@adsc.com.sg, eleyans@nus.edu.sg,

moulin@ifp.uiuc.edu, qi.tian@utsa.edu

**Abstract.** In fine-grained action (object manipulation) recognition, it is important to encode object semantic (contextual) information, i.e., which object is being manipulated and how it is being operated. However, previous methods for action recognition often represent the semantic information in a global and coarse way and therefore cannot cope with fine-grained actions. In this work, we propose a representation and classification pipeline which seamlessly incorporates localized semantic information into every processing step for fine-grained action recognition. In the feature extraction stage, we explore the geometric information between local motion features and the surrounding objects. In the feature encoding stage, we develop a semantic-grouped locality-constrained linear coding (SG-LLC) method that captures the joint distributions between motion and object-in-use information. Finally, we propose a semantic-aware multiple kernel learning framework (SA-MKL) by utilizing the empirical joint distribution between action and object type for more discriminative action classification. Extensive experiments are performed on the large-scale and difficult fine-grained MPII cooking action dataset. The results show that by effectively accumulating localized semantic information into the action representation and classification pipeline, we significantly improve the fine-grained action classification performance over the existing methods.

1

\*\*\*\*\*

## Robust Scene Text Detection with Convolution Neural Network Induced MSER Trees

Weilin Huang<sup>1,2</sup>, Yu Qiao<sup>1</sup>, and Xiaoou Tang<sup>2,1</sup>

<sup>1</sup>Shenzhen Key Lab of Comp. Vis and Pat. Rec.,

Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, China

<sup>2</sup>Department of Information Engineering,  
The Chinese University of Hong Kong, China

**Abstract.** Maximally Stable Extremal Regions (MSERs) have achieved great success in scene text detection. However, this low-level pixel operation inherently limits its capability for handling complex text information efficiently (e. g. connections between text or background components), leading to the difficulty in distinguishing texts from background components. In this paper, we propose a novel framework to tackle this problem by leveraging the high capability of convolutional neural network (CNN). In contrast to recent methods using a set of low-level heuristic features, the CNN network is capable of learning high-level features to robustly identify text components from text-like outliers (e.g. bikes, windows, or leaves). Our approach takes advantages of both MSERs and sliding-window based methods. The MSERs operator dramatically reduces the number of windows scanned and enhances detection of the low-quality texts. While the sliding-window with CNN is applied to correctly separate the connections of multiple characters in components. The proposed system achieved strong robustness against a number of extreme text variations and serious real-world problems. It was evaluated on the ICDAR 2011 benchmark dataset, and achieved over 78% in F-measure, which is significantly higher than previous methods.

**Keywords:** Maximally Stable Extremal Regions (MSERs), convolutional neural network (CNN), text-like outliers, sliding-window.

1

\*\*\*\*\*

Deep Features for Text Spotting

Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman  
Visual Geometry Group, Department of Engineering Science,  
University of Oxford, UK

**Abstract.** The goal of this work is text spotting in natural images. This is divided into two sequential tasks: detecting words regions in the image, and recognizing the words within these regions. We make the following contributions: First, we develop a Convolutional Neural Network (CNN) classifier that can be used for both tasks. The CNN has a novel architecture that enables efficient feature sharing (by using a number of layers in common) for text detection, character case-sensitive and insensitive classification, and bigram classification. It exceeds the state-of-the-art performance for all of these. Second, we make a number of technical changes over the traditional CNN architectures, including no downsampling for a per-pixel sliding window, and multi-mode learning with a mixture of linear models (maxout). Third, we have a method of automated data mining of Flickr, that generates word and character level annotations. Finally, these components are used together to form an end-to-end, state-of-the-art text spotting system. We evaluate the text-spotting system on two standard benchmarks, the ICDAR Robust Reading data set and the Street View Text data set, and demonstrate improvements over the state-of-the-art on multiple measures.

1

\*\*\*\*\*

Improving Image-Sentence Embeddings Using  
Large Weakly Annotated Photo Collections

Yunchao Gong<sup>1</sup>, Liwei Wang<sup>2</sup>, Michal Hodosh<sup>2</sup>, Julia Hockenmaier<sup>2</sup>,  
and Svetlana Lazebnik<sup>2</sup>

<sup>1</sup>University of North Carolina at Chapel Hill, USA  
yunchao@cs.unc.edu

<sup>2</sup>University of Illinois at Urbana-Champaign, USA  
{lwang97, mhodosh2, juliahmr, slazebni}@illinois.edu

**Abstract.** This paper studies the problem of associating images with descriptive sentences by embedding them in a common latent space. We are interested in learning such embeddings from hundreds of thousands

or millions of examples. Unfortunately, it is prohibitively expensive to fully annotate this many training images with ground-truth sentences. Instead, we ask whether we can learn better image-sentence embeddings by augmenting small fully annotated training sets with millions of images that have weak and noisy annotations (titles, tags, or descriptions). After investigating several state-of-the-art scalable embedding methods, we introduce a new algorithm called Stacked Auxiliary Embedding that can successfully transfer knowledge from millions of weakly annotated images to improve the accuracy of retrieval-based image description.

1

\*\*\*\*\*

## Strengthening the Effectiveness of Pedestrian Detection with Spatially Pooled Features

Sakrapee Paisitkriangkrai, Chunhua Shen\*, and Anton van den Hengel  
The University of Adelaide, Australia  
chunhua.shen@adelaide.edu.au

Abstract. We propose a simple yet effective approach to the problem of pedestrian detection which outperforms the current state-of-the-art. Our new features are built on the basis of low-level visual features and spatial pooling. Incorporating spatial pooling improves the translational invariance and thus the robustness of the detection process. We then directly optimise the partial area under the ROC curve (pAUC) measure, which concentrates detection performance in the range of most practical importance. The combination of these factors leads to a pedestrian detector which outperforms all competitors on all of the standard benchmark datasets. We advance state-of-the-art results by lowering the average miss rate from 13% to 11% on the INRIA benchmark, 41% to 37% on the ETH benchmark, 51% to 42% on the TUD-Brussels benchmark and 36% to 29% on the Caltech-USA benchmark.

1

\*\*\*\*\*

## Selecting Influential Examples: Active Learning with Expected Model Output Changes

Alexander Freytag, Erik Rodner, and Joachim Denzler  
Computer Vision Group, Friedrich Schiller University Jena, Germany  
{firstname.lastname}@uni-jena.de  
<http://www.inf-cv.uni-jena.de>

Abstract. In this paper, we introduce a new general strategy for active learning.

The key idea of our approach is to measure the expected change of model outputs, a concept that generalizes previous methods based on expected model change and incorporates the underlying data distribution. For each example of an unlabeled set, the expected change of model predictions is calculated and marginalized over

the unknown label. This results in a score for each unlabeled example that can be

used for active learning with a broad range of models and learning algorithms. In

particular, we show how to derive very efficient active learning methods for Gaussian

process regression, which implement this general strategy, and link them to previous methods. We analyze our algorithms and compare them to a broad range of previous active learning strategies in experiments showing that they outperform state-of-the-art on well-established benchmark datasets in the area of visual

object recognition.

Keywords: active learning, Gaussian processes, visual recognition, exploration-exploitation trade-off.

1

\*\*\*\*\*

## Efficient Sparsity Estimation via Marginal-Lasso Coding

Tzu-Yi Hung<sup>1</sup>, Jiwen Lu<sup>2</sup>, Yap-Peng Tan<sup>1</sup>, and Shenghua Gao<sup>3</sup>

<sup>1</sup>School of Electrical and Electronic Engineering,  
Nanyang Technological University, Singapore

<sup>2</sup>Advanced Digital Sciences Center, Singapore

<sup>3</sup>ShanghaiTech University, Shanghai, China

**Abstract.** This paper presents a generic optimization framework for efficient feature quantization using sparse coding which can be applied to many computer vision tasks. While there are many works working on sparse coding and dictionary learning, none of them has exploited the advantages of the marginal regression and the lasso simultaneously to provide more efficient and effective solutions. In our work, we provide such an approach with a theoretical support. Therefore, the

computational complexity of the proposed method can be two orders faster than that of the lasso with sacrificing the inevitable quantization error. On the other hand, the proposed method is more robust than the conventional marginal regression based methods. We also provide an adaptive regularization parameter selection scheme and a dictionary learning method incorporated with the proposed sparsity estimation algorithm. Experimental results and detailed model analysis are presented to demonstrate the efficacy of our proposed methods.

**Keywords:** Sparsity estimation, marginal regression, sparse coding, lasso, dictionary learning, adaptive regularization parameter.

1

\*\*\*\*\*

## Continuous Conditional Neural Fields for Structured Regression

Tadas Baltrušaitis<sup>1</sup>, Peter Robinson<sup>1</sup>, and Louis-Philippe Morency<sup>2</sup>

<sup>1</sup>Computer Laboratory, University of Cambridge, UK

{tadas.baltrušaitis, peter.robinson}@cl.cam.ac.uk

<sup>2</sup>Institute for Creative Technologies, University of Southern California, CA  
morency@ict.usc.edu

**Abstract.** An increasing number of computer vision and pattern recognition problems require structured regression techniques. Problems like human pose estimation, unsegmented action recognition, emotion prediction and facial landmark detection have temporal or spatial output dependencies that regular regression techniques do not capture. In this paper we present continuous conditional neural fields (CCNF) - a novel structured regression model that can learn non-linear input-output dependencies, and model temporal and spatial output relationships of varying length sequences. We propose two instances of our CCNF framework:

Chain-CCNF for time series modelling, and Grid-CCNF for spatial relationship modelling. We evaluate our model on five public datasets spanning three different regression problems: facial landmark detection in the wild, emotion prediction in music and facial action unit recognition. Our CCNF model demonstrates state-of-the-art performance on all of the datasets used.

**Keywords:** Structured regression, Landmark detection, Face tracking.

1

\*\*\*\*\*

## Learning to Rank Using High-Order Information

Puneet Kumar Dokania<sup>1</sup>, Aseem Behl<sup>2</sup>, Chao V. Jawahar<sup>2</sup>, and M. Pawan Kumar<sup>1</sup>

<sup>1</sup>Ecole Centrale de Paris

INRIA Saclay, France

<sup>2</sup>IIIT Hyderabad, India

**Abstract.** The problem of ranking a set of visual samples according to their relevance to a query plays an important role in computer vision. The traditional approach for ranking is to train a binary classifier such as a support vector machine (svm). Binary classifiers suffer from two main deficiencies: (i) they do not optimize a ranking-based loss function, for example, the average precision (ap) loss; and (ii) they cannot incorporate



high-order information such as the a p r i o r i correlation between the relevance of two visual samples (for example, two persons in the same image tend to perform the same action). We propose two novel learning formulations that allow us to incorporate high-order information for ranking.

The first framework, called high-order binary svm(hob-svm), allows for a structured input. The parameters of hob-svm are learned by minimizing a convex upper bound on a surrogate 0-1 loss function. In order to obtain the ranking of the samples that form the structured input, hob-svm sorts the samples according to their max-marginals. The second framework, called high-order average precision svm(hoap-svm), also allows for a structured input and uses the same ranking criterion. However, in contrast to hob-svm, the parameters of hoap-svm are learned by minimizing a difference-of-convex upper bound on the a p l o s s. Using a standard, publicly available dataset for the challenging problem of action classification, we show that both hob-svm and hoap-svm outperform the baselines that ignore high-order information.

1

\*\*\*\*\*

Support Vector Guided Dictionary Learning

Sijia Cai<sup>1,3</sup>, Wangmeng Zuo<sup>2</sup>, L e i Z h a n g<sup>3,4</sup>, Xiangchu Feng<sup>4</sup>, and Ping Wang<sup>1</sup>

<sup>1</sup>School of Science, Tianjin University, China

<sup>2</sup>School of Computer Science and Technology, Harbin Institute of Technology, China

<sup>3</sup>Dept. of Computing, The Hong Kong Polytechnic University, China

<sup>4</sup>Dept. of Applied Mathematics, Xidian University, China

cssjcai@gmail.com, cslzhang@comp.polyu.edu.hk

**Abstract.** Discriminative dictionary learning aims to learn a dictionary from training samples to enhance the discriminative capability of their coding vectors. Several discrimination terms have been proposed by assessing the prediction loss (e.g., logistic regression) or class separation criterion (e.g., Fisher discrimination criterion) on the coding vectors. In this paper, we provide a new insight on discriminative dictionary learning. Specifically, we formulate the discrimination term as the weighted summation of the squared distances between all pairs of coding vectors. The discrimination term in the state-of-the-art Fisher discrimination dictionary learning (FDDL) method can be explained as a special case of our model, where the weights are simply determined by the numbers of samples of each class. We then propose a parameterization method to adaptively determine the weight of each coding vector pair, which leads to a support vector guided dictionary learning (SVGDL) model. Compared with FDDL, SVGDL can adaptively assign different weights to different pairs of coding vectors. More importantly, SVGDL automatically selects only a few critical pairs to assign non-zero weights, resulting in better generalization ability for pattern recognition tasks. The experimental results on a series of benchmark databases show that SVGDL outperforms many state-of-the-art discriminative dictionary learning methods.

**Keywords:** Dictionary learning, support vector machine, sparse representation, Fisher discrimination.

1

\*\*\*\*\*

Video Object Discovery and Co-segmentation  
with Extremely Weak Supervision

Le Wang<sup>1</sup>, G a n g H u a<sup>2</sup>, Rahul Sukthankar<sup>3</sup>, Jianru Xue<sup>1</sup>, and Nanning Zheng<sup>1</sup>

<sup>1</sup>Xi'an Jiaotong University, China

<sup>2</sup>Stevens Institute of Technology, USA

<sup>3</sup>Google Research, USA

**Abstract.** Video object co-segmentation refers to the problem of simultaneously segmenting a common category of objects from multiple videos. Most existing video co-segmentation methods assume that all frames from all videos contain the target objects. Unfortunately, this assumption is rarely true in practice, p

ar-  
 ticularly for large video sets, and existing methods perform poorly when the  
 assumption is violated. Hence, any practical video object co-segmentation al-  
 gorithm needs to identify the relevant frames containing the target object from  
 all videos, and then co-segment the object only from these relevant frames. We  
 present a spatiotemporal energy minimization formulation for simultaneous video  
 object discovery and co-segmentation across multiple videos. Our formulation in-  
 corporates a spatiotemporal auto-context model, which is combined with appear-  
 ance modeling for superpixel labeling. The superpixel-level labels are propagate  
 d  
 to the frame level through a multiple instance boosting algorithm with spatial  
 rea-  
 soning (Spatial-MILBoosting), based on which frames containing the video ob-  
 ject are identi-  
 fied. Our method only needs to be bootstrapped with the frame-level  
 labels for a few video frames ( e.g., usually 1 to 3) to indicate if they contain the  
 target objects or not. Experiments on three datasets validate the ef-  
 ficacy of our  
 proposed method, which compares favorably with the state-of-the-art.  
 Keywords: video object discovery, video object co-segmentation, spatiotempo-  
 ral auto-context model, Spatial-MILBoosting.

1  
 \*\*\*\*\*

Supervoxel-Consistent  
 Foreground Propagation in Video  
 Suyog Dutt Jain and Kristen Grauman  
 University of Texas at Austin, USA

Abstract. A major challenge in video segmentation is that the fore-  
 ground object may move quickly in the scene at the same time its appearance and  
 shape evolves over time. While pairwise potentials used  
 in graph-based algorithms help smooth labels between neighboring (su-  
 per)pixels in space and time, they offer only a myopic view of consistency and c  
 an be misled by inter-frame optical flow errors. We propose  
 a hierarchical supervoxel label consistency potential for semi-supervised  
 foreground segmentation. Given an initial frame with manual annotation for the  
 foreground object, our approach propagates the foreground  
 region through time, leveraging bottom-up supervoxels to guide its es-  
 timates towards long-range coherent regions. We validate our approach on three ch  
 allenging datasets and achieve state-of-the-art results.

1  
 \*\*\*\*\*

Clustering with Hypergraphs:  
 The Case for Large Hyperedges  
 Pulak Purkait<sup>1</sup>, Tat-Jung Chin<sup>1</sup>, Hanno Ackermann<sup>2</sup>, and David Suter<sup>1</sup>  
<sup>1</sup>The University of Adelaide, Australia  
<sup>2</sup>Leibniz Universität Hannover, Germany

Abstract. The extension of conventional clustering to hypergraph clus-  
 tering, which involves higher order similarities instead of pairwise simi-  
 larities, is increasingly gaining attention in computer vision. This is due  
 to the fact that many grouping problems require an affinity measure that must invol  
 ve a subset of data of size more than two, i.e., a hyperedge. Almost all previous works, however, have considered the smallest possible  
 hyperedge size, due to a lack of study into the potential benefits of large  
 hyperedges and effective algorithms to generate them. In this paper, we show that  
 large hyperedges are better from both theoretical and empirical  
 standpoints. We then propose a novel guided sampling strategy for  
 large hyperedges, based on the concept of random cluster models. Our  
 method can generate pure large hyperedges that significantly improve  
 grouping accuracy without exponential increases in sampling costs. In  
 the important applications of face clustering and motion segmentation, our method

demonstrates substantially better accuracy and efficiency.  
Keywords: Hypergraph clustering, model fitting, guided sampling.

1

\*\*\*\*\*

Person Re-identification by Video Ranking

Taiqing Wang<sup>1</sup>, Shaogang Gong<sup>2</sup>, Xiatian Zhu<sup>2</sup>, and Shengjin Wang<sup>1</sup>

<sup>1</sup>Dept. of Electronic Engineering, Tsinghua University, China

<sup>2</sup>School of EECS, Queen Mary University of London, UK

Abstract. Current person re-identification (re-id) methods typically rely on single-frame imagery features, and ignore space-time information from image sequences. Single-frame (single-shot) visual appearance matching is inherently limited for person re-id in public spaces due to visual ambiguity arising from non-overlapping camera views where viewpoint and lighting changes can cause significant appearance variation. In this work, we present a novel model to automatically select the most discriminative video fragments from noisy image sequences of people where more reliable space-time features can be extracted, whilst simultaneously to learn a video ranking function for person re-id. Also, we introduce a new image sequence re-id dataset (iLIDS-VID) based on the iLIDS-MCT benchmark data. Using the iLIDS-VID and PRID2011 sequence re-id datasets, we extensively conducted comparative evaluations to demonstrate the advantages of the proposed model over contemporary gait recognition, holistic image sequence matching and state-of-the-art single-shot/multi-shot based re-id methods.

1

\*\*\*\*\*

Bayesian Nonparametric

Intrinsic Image Decomposition

Jason Chang, Randi Cabezas, and John W. Fisher III

CSAIL, MIT, USA

Abstract. We present a generative, probabilistic model that decomposes an image into reflectance and shading components. The proposed approach uses a Dirichlet process Gaussian mixture model where the mean parameters evolve jointly according to a Gaussian process. In contrast to prior methods, we eliminate the Retinex term and adopt more general smoothness assumptions for the shading image. Markov chain Monte Carlo sampling techniques are used for inference, yielding state-of-the-art results on the MIT Intrinsic Image Dataset.

Keywords: Intrinsic images, Dirichlet process, Gaussian process, MCMC.

1

\*\*\*\*\*

Face Detection without Bells and Whistles

Markus Mathias<sup>1</sup>, Rodrigo Benenson<sup>2</sup>, Marco Pedersoli<sup>1</sup>, and Luc Van Gool<sup>1,3</sup>

<sup>1</sup>ESAT-PSI/VISICS, iMinds, KU Leuven, Belgium

<sup>2</sup>MPI Informatics, Saarbrücken, Germany

<sup>3</sup>3D-ITET/CVL, ETH Zürich, Switzerland

Abstract. Face detection is a mature problem in computer vision. While diverse high performing face detectors have been proposed in the past, we present two surprising new top performance results. First, we show that a properly trained vanilla DPM reaches top performance, improving over commercial and research systems. Second, we show that a detector based on rigid templates - similar in structure to the Viola&Jones detector - can reach similar top performance on this task. Importantly, we discuss issues with existing evaluation benchmark and propose an improved procedure.

Fig. 1. Our proposed HeadHunter detector at the Oscars. Can you spot the one false

positive, and one false negatives? (hint: first rows).

1

\*\*\*\*\*

## On Image Contours of Projective Shapes

Jean Ponce<sup>1</sup>, and Martial Hebert<sup>2</sup>

<sup>1</sup>Department of Computer Science,  
Ecole Normale Supérieure, France

<sup>2</sup>Robotics Institute,  
Carnegie-Mellon University, USA

**Abstract.** This paper revisits classical properties of the outlines of solid shapes

bounded by smooth surfaces, and shows that they can be established in a purely projective setting, without appealing to Euclidean measurements such as normals or curvatures. In particular, we give new synthetic proofs of Koenderink's famous

theorem on convexities and concavities of the image contour, and of the fact that

the rim turns in the same direction as the viewpoint in the tangent plane at a convex point, and in the opposite direction at a hyperbolic point. This suggests that projective geometry should not be viewed merely as an analytical device for linearizing calculations (its main role in structure from motion), but as the

proper framework for studying the relation between solid shape and its perspective projections. Unlike previous work in this area, the proposed approach does not require an oriented setting, nor does it rely on any choice of coordinate system or analytical considerations.

1

\*\*\*\*\*

## Programmable Automotive Headlights

Robert Tamburo<sup>1</sup>, Eriko Nurvitadhi<sup>2</sup>, Abhishek Chugh<sup>1</sup>, Mei Chen<sup>2</sup>,

Anthony Rowel, Takeo Kanade<sup>1</sup>, and Srinivasa G. Narasimhan<sup>1</sup>

<sup>1</sup>Carnegie Mellon University, Pittsburgh, PA, USA

<sup>2</sup>Intel Labs, Pittsburgh, PA, USA

**Abstract.** The primary goal of an automotive headlight is to improve safety in low light and poor weather conditions. But, despite decades of innovation on light sources, more than half of accidents occur at night

even with less traffic on the road. Recent developments in adaptive lighting have addressed some limitations of standard headlights, however, they have limited flexibility - switching between high and low beams,

turning off beams toward the opposing lane, or rotating the beam as the vehicle turns - and are not designed for all driving environments. This paper introduces an ultra-low latency reactive visual system that

can sense, react, and adapt quickly to any environment while moving at highway speeds. Our single hardware design can be programmed to perform a variety of tasks. Anti-glare high beams, improved driver visibility

during snowstorms, increased contrast of lanes, markings, and sidewalks, and early visual warning of obstacles are demonstrated.

**Keywords:** Adaptive headlights, reactive visual system, computational illumination.

1

\*\*\*\*\*

## ROCHADE: Robust Checkerboard Advanced

Detection for Camera Calibration

Simon Placht<sup>1,2</sup>, Peter F. F. Sattler<sup>1,2</sup>, Etienne Assoumou Mengue<sup>2</sup>,

Hannes Hofmann<sup>1</sup>, Christian Schaller<sup>1</sup>, Michael Balda<sup>1</sup>,

and Elli Angelopoulou<sup>2</sup>

<sup>1</sup>Metrilus GmbH, Erlangen, Germany

<sup>2</sup>Pattern Recognition Lab, University of Erlangen, Nuremberg, Germany

**Abstract.** We present a new checkerboard detection algorithm which is able to detect checkerboards at extreme poses, or checkerboards which are highly distorted due to lens distortion even on low-resolution images.

On the detected pattern we apply a surface fitting based subpixel re-

Refinement specifically tailored for checkerboard X-junctions. Finally, we investigate how the accuracy of a checkerboard detector affects the overall calibration result in multi-camera setups. The proposed method is evaluated on real images captured with different camera models to show its wide applicability. Quantitative comparisons to OpenCV's checkerboard detector show that the proposed method detects up to 80% more checkerboards and detects corner points more accurately, even under strong perspective distortion as often present in wide baseline stereo setups.

Keywords: Checkerboard Detection, Saddle-Based Subpixel Refinement, Multi Camera Calibration, Low Resolution Sensors, Lens Distortion.

1

\*\*\*\*\*

#### Correcting for Duplicate Scene Structure in Sparse 3D Reconstruction

Jared Heinly, Enrique Dunn, and Jan-Michael Frahm

The University of North Carolina at Chapel Hill, USA

Abstract. Structure from motion (SfM) is a common technique to recover 3D geometry and camera poses from sets of images of a common scene. In many urban environments, however, there are symmetric, repetitive, or duplicate structures that pose challenges for SfM pipelines. The result of these ambiguous structures is incorrectly placed cameras and points within the reconstruction. In this paper, we present a post-processing method that can not only detect these errors, but successfully resolve them. Our novel approach proposes the strong and informative measure of conflicting observations, and we demonstrate that it is robust to a large variety of scenes.

Keywords: Structure from motion, duplicate structure disambiguation.

1

\*\*\*\*\*

#### Total Moving Face Reconstruction

Supasorn Suwajanakorn, Ira Kemelmacher-Shlizerman, and Steven M. Seitz

University of Washington, USA

Fig. 1. Given a YouTube video of a person's face our method estimates high detail geometry (full 3D flow and pose) in each video frame completely automatically

Abstract. We present an approach that takes a single video of a person's face and reconstructs a high detail 3D shape for each video frame.

We target videos taken under uncontrolled and uncalibrated imaging conditions, such as YouTube videos of celebrities. In the heart of this work is a new dense 3D flow estimation method coupled with shape from shading. Unlike related works we do not assume availability of a blend shape model, nor require the person to participate in a training/capturing process. Instead we leverage the large amounts of photos that are available per individual in personal or internet photo collections. We show results for a variety of video sequences that include various lighting conditions, head poses, and facial expressions.

Keywords: 3D reconstruction, faces, non-rigid reconstruction.

1

\*\*\*\*\*

#### Automatic Single-View Calibration and Rectification from Parallel Planar Curves

Eduardo R. Corral-Soto and James H. Elder

Centre for Vision Research, York University, Toronto, Canada

Abstract. Typical methods for camera calibration and image rectification from a single view assume the existence of straight parallel lines from which vanishing points can be computed, or orthogonal structure known to exist in the scene. However, there are practical situations where these assumptions do not apply. Moreover, from a single family of parallel lines on the ground plane there is insufficient information to recover a complete rectification. Here we study a generalization of these meth-

ods to scenes known to contain parallel curves. Our method is based on establishing an association between pairs of corresponding points lying on the image projection of these curves. We show how this method can be used to compute a least-squares estimate of the focal length and the camera pose from the tangent lines of the associated points, allowing complete rectification of the image. We evaluate the method on highway and sports track imagery, and demonstrate its accuracy relative to a state-of-the-art vanishing point method.

Keywords: camera calibration, projective rectification, contour grouping, traffic surveillance.

1

\*\*\*\*\*

On Sampling Focal Length Values

to Solve the Absolute Pose Problem

Torsten Sattler<sup>1</sup>, Chris Sweeney<sup>2,✉</sup>, and Marc Pollefeys<sup>1</sup>

<sup>1</sup>Department of Computer Science, ETH Zürich, Zürich, Switzerland

<sup>2</sup>University of California Santa Barbara, Santa Barbara, USA

Abstract. Estimating the absolute pose of a camera relative to a 3D representation of a scene is a fundamental step in many geometric Computer Vision applications. When the camera is calibrated, the pose can be computed very efficiently. If the calibration is unknown, the problem becomes much harder, resulting in slower solvers or solvers requiring more samples and thus significantly longer run-times for RANSAC. In this paper, we challenge the notion that using minimal solvers is always optimal and propose to compute the pose for a camera with unknown focal length by randomly sampling a focal length value and using an efficient pose solver for the now calibrated camera. Our main contribution is a novel sampling scheme that enables us to guide the sampling process towards promising focal length values and avoids considering all possible values once a good pose is found. The resulting RANSAC variant is significantly faster than current state-of-the-art pose solvers, especially for low inlier ratios, while achieving a similar or better pose accuracy.

Keywords: RANSAC, n-point-pose (P nP), camera pose estimation.

1

\*\*\*\*\*