

Lie Bodies: A Manifold Representation
of 3D Human Shape

Oren Freifeld¹ and Michael J. Black²

¹Division of Applied Mathematics, Brown University, Providence, RI 02912, USA

²Max Planck Institute for Intelligent Systems, 72076 Tübingen, Germany

freifeld@dam.brown.edu, black@is.mpg.de

Abstract. Three-dimensional object shape is commonly represented in terms of deformations of a triangular mesh from an exemplar shape. Existing models, however, are based on a Euclidean representation of shape deformations. In contrast, we argue that shape has a manifold structure:

For example, summing the shape deformations for two people does not necessarily yield a deformation corresponding to a valid human shape, nor does the Euclidean difference of these two deformations provide a

meaningful measure of shape dissimilarity. Consequently, we define a novel manifold for shape representation, with emphasis on body shapes, using a new Lie group of deformations. This has several advantages. First

we define triangle deformations exactly, removing non-physical deformations and redundant degrees of freedom common to previous methods. Second, the Riemannian structure of Lie Bodies enables a more mean-

ingful definition of body shape similarity by measuring distance between

bodies on the manifold of body shape deformations. Third, the group

structure allows the valid composition of deformations. This is important

for models that factor body shape deformations into multiple causes or

represent shape as a linear combination of basis shapes. Finally, body

shape variation is modeled using statistics on manifolds. Instead of mod-

eling Euclidean shape variation with Principal Component Analysis we

capture shape variation on the manifold using Principal Geodesic Analysis. Our

experiments show consistent visual and quantitative advantages of Lie Bodies over traditional Euclidean models of shape deformation

and our representation can be easily incorporated into existing methods.

Keywords: Shape deformation, Lie group, Statistics on manifolds.

1

Worldwide Pose Estimation Using 3D Point Clouds

Yunpeng Li¹, Noah Snavely², Dan Huttenlocher², and Pascal Fua¹

¹EPFL

{yunpeng.li, pascal.fua}@epfl.ch

²Cornell University

{snavely, dph}@cs.cornell.edu

Abstract. We address the problem of determining where a photo was taken by estimating a full 6-DOF-plus-intrinsic camera pose with respect to a large geo-registered 3D point cloud, bringing together research on image localization, land-

mark recognition, and 3D pose estimation. Our method scales to datasets with hundreds of thousands of images and tens of millions of 3D points through the use of two new techniques: a co-occurrence prior for RANSAC and bidirectional matching of image features with 3D points. We evaluate our method on several large data sets, and show state-of-the-art results on landmark recognition as well

as the ability to locate camera as to within meters, requiring only seconds per query.

1

Improved Reconstruction of Deforming Surfaces

by Cancelling Ambient Occlusion

Thabo Beeler^{1,2}, Derek Bradley¹, Henning Zimmer², and Markus Gross^{1,2}

¹Disney Research Zurich

{derek.bradley, dbeeler}@disneyresearch.com

²ETH Zurich

{hzimmer, grossm}@inf.ethz.ch

Abstract. We present a general technique for improving space-time reconstruction of deforming surfaces, which are captured in a video-based reconstruction scenario under uniform illumination. Our approach simultaneously improves both the acquired shape as well as the tracked motion of the deforming surface. The method is based on factoring out surface shading, computed by a fast approximation to global illumination called ambient occlusion. This allows us to improve the performance of optical flow tracking that mainly relies on constancy of image features, such as intensity. While cancelling the local shading, we also optimize the surface shape to minimize the residual between the ambient occlusion of the 3D geometry and that of the image, yielding more accurate surface details in the reconstruction. Our enhancement is independent of the actual space-time reconstruction algorithm. We experimentally measure the quantitative improvements produced by our algorithm using a synthetic example of deforming skin, where ground truth shape and motion is available. We further demonstrate our enhancement on a real-world sequence of human face reconstruction.

1

On the Statistical Determination of Optimal Camera Configurations in Large Scale Surveillance Networks

Junbin Liu¹, Clinton Fookes¹, Tim Wark², and Sridha Sridharan¹

¹Image & Video Research Laboratory, Queensland University of Technology,

2 George Street, Brisbane, QLD 4000, Australia

junbin.liu@gmail.com, {c.fookes,s.sridhara}@qut.edu.au

²CSIRO ICT Centre,

1 Technology Court, Pullenvale, QLD 4069, Australia

tim.wark@csiro.au

Abstract. The selection of optimal camera configurations (camera locations, orientations etc.) for multi-camera networks remains an unsolved problem. Previous approaches largely focus on proposing various objective functions to achieve different tasks. Most of them, however, do not generalize well to large scale networks.

To tackle this, we introduce a statistical formulation of the optimal selection

of camera configurations as well as propose a Trans-Dimensional Simulated Annealing (TDSA) algorithm to effectively solve the problem. We compare our approach with a state-of-the-art method based on Binary Integer Programming (BIP) and show that our approach offers similar performance on small scale problems. However, we also demonstrate the capability of our approach in dealing with large scale problems and show that our approach produces better results than 2 alternative

heuristics designed to deal with the scalability issue of BIP.

Keywords: Camera placement, optimization, resolvable jump Markov chain Monte Carlo, simulated annealing.

1

The Scale of Geometric Texture

Geoffrey Oxholm, Prabin Bariya, and Ko Nishino

Department of Computer Science

Drexel University, Philadelphia, PA 19104, USA

{gao25,pb335,kon}@drexel.edu

Abstract. The most defining characteristic of texture is its underlying geometry. Although the appearance of texture is as dynamic as its illumination and viewing conditions, its geometry remains constant. In this work, we study the fundamental characteristic properties of texture geometry—self similarity and scale variability—and exploit them to perform surface normal estimation, and geometric texture classification. Textures, whether they are regular or stochastic, exhibit

some form of repetition in their underlying geometry. We use this property to

derive a photometric stereo method uniquely tailored to utilize the redundancy in geometric texture. Using basic observations about the scale variability of texture geometry, we derive a compact, rotation invariant, scale-space representation of geometric texture. To evaluate this representation we introduce an extensive new texture database that contains multiple distances as well as in-plane and out-of plane rotations. The high accuracy of the classification results indicate the descriptive yet compact nature of our texture representation, and demonstrates the importance of geometric texture analysis, pointing the way towards improvements in appearance modeling and synthesis.

1

Efficient Articulated Trajectory Reconstruction

Using Dynamic Programming and Filters

Jack Valmadre^{1,3}, Yingying Zhu^{2,3}, Sridha Sridharan¹, and Simon Lucey^{1,3}

¹Queensland University of Technology, Australia

²University of Queensland, Australia

³Commonwealth Scientific and Industrial Research Organisation, Australia

{jack.valmadre,yingying.zhu,simon.lucey}@csiro.au

Abstract. This paper considers the problem of reconstructing the motion of a 3D articulated tree from 2D point correspondences subject to some temporal prior. Hitherto, smooth motion has been encouraged using a trajectory basis, yielding a hard combinatorial problem with time complexity growing exponentially in the number of frames. Branch and bound strategies have previously attempted to curb this complexity whilst maintaining global optimality. However, they provide no guarantee of being more efficient than exhaustive search. Inspired by recent work which reconstructs general trajectories using compact high-pass filters, we develop a dynamic programming approach which scales linearly in the number of frames, leveraging the intrinsically local nature of filter interactions. Extension to affine projection enables reconstruction without estimating cameras.

1

Object Co-detection

Sid Yingze Bao, Yu Xiang, and Silvio Savarese

University of Michigan at Ann Arbor, USA

{yingze,yuxiang,silvio}@eecs.umich.edu

Abstract. In this paper we introduce a new problem which we call object co-detection. Given a set of images with objects observed from two or multiple images, the goal of co-detection is to detect the objects, establish the identity of individual object instance, as well as estimate the viewpoint transformation of corresponding object instances. In designing a co-detector, we follow the intuition that an object has consistent appearance when observed from the same or different viewpoints. By modeling an object using state-of-the-art part-based representations such as [1,2], we measure appearance consistency between objects by comparing part appearance and geometry across images. This allows to effectively account for object self-occlusions and viewpoint transformations. Extensive experimental evaluation indicates that our co-detector obtains more accurate detection results than if objects were to be detected from each image individually. Moreover, we demonstrate the relevance of our co-detection scheme to other recognition problems such as single instance object recognition, wide-baseline matching, and image query.

1

Morphable Displacement Field Based Image

Matching for Face Recognition across Pose

Shaoxin Li^{1,2}, Xin Liu^{1,2}, Xiujuan Chai¹, Haihong Zhang³,
Shihong Lao³, and Shiguang Shan¹

¹Key Lab of Intelligent Information Processing of Chinese Academy of Sciences
(CAS), Institute of Computing Technology, CAS, Beijing, 100190, China

²Graduate University of Chinese Academy of Sciences, Beijing 100049, China

³Omron Social Solutions Co., LTD., Kyoto, Japan

{shaoxin.li,xiujuan.chai,xin.liu,shiguang.shan}@vipl.ict.ac.cn,

lao@ari.ncl.omron.co.jp, angelazhang@ssb.kusatsu.omron.co.jp

Abstract. Fully automatic Face Recognition Across Pose (FRAP) is one of the most desirable techniques, however, also one of the most challenging tasks in face recognition. Matching a pair of face images in different poses can be converted into matching their pixels corresponding to the same semantic facial point. Following this idea, given two images from different poses, we propose a novel method, named Morphable Displacement Field (MDF), to match with the virtual view under the pose. By formulating MDF as a convex combination of a number of template displacement fields generated from a 3D face database, our model satisfies both global conformity and local consistency. We further present an approximate but effective solution of the proposed MDF model, named implicit Morphable Displacement Field (iMDF), which synthesizes virtual view implicitly via an MDF by minimizing matching residual. This formulation not only avoids intractable optimization of the high-dimensional displacement field but also facilitates a constrained quadratic optimization. The proposed method can work well even when only 2 facial landmarks are labeled, which makes it especially suitable for fully automatic FRAP system. Extensive evaluations on FERET, PIE and Multi-PIE databases show considerable improvement over state-of-the-art FRAP algorithms in both semi-automatic and fully automatic evaluation protocols.

1

Combining Per-frame and Per-track Cues

for Multi-person Action Recognition

Sameh Khamis, Vlad I. Morariu, and Larry S. Davis

University of Maryland, College Park

{sameh,morariu,lsd}@umiacs.umd.edu

Abstract. We propose a model to combine per-frame and per-track cues for action recognition. With multiple targets in a scene, our model simultaneously captures the natural harmony of an individual's action in a scene and the flow of actions of an individual in a video sequence, inferring valid tracks in the process. Our motivation is based on the unlikely discordance of an action in a structured scene, both at the track level and the frame level (e.g., a person dancing in a crowd of joggers). While we can utilize sampling approaches for inference in our model, we instead devise a global inference algorithm by decomposing the problem and solving the subproblems exactly and efficiently, recovering a globally optimal joint solution in several cases. Finally, we improve on the state-of-the-art action recognition results for two publicly available datasets.

1

Joint Image and Word Sense Discrimination

for Image Retrieval

Aurelien Lucchi^{1,2} and Jason Weston¹

¹Google, New York, USA

²EPFL, Lausanne, Switzerland

Abstract. We study the task of learning to rank images given a text query, a problem that is complicated by the issue of multiple senses. That is, the senses of interest are typically the visually distinct concepts that a user wishes to retrieve. In this paper, we propose to learn a ranking function that

optimizes the ranking cost of interest and simultaneously discovers the disambiguated senses of the query that are optimal for the supervised task. Note that no supervised information is given about the senses. Experiments performed on web images and the ImageNet dataset show that using our approach leads to a clear gain in performance.

1

Script Data for Attribute-Based Recognition of Composite Activities

Marcus Rohrbach¹, Michaela Regneri², Mykhaylo Andriluka¹,
Sikandar Amin^{1,3}, Manfred Pinkal², and Bernt Schiele¹

¹Max Planck Institute for Informatics, Saarbrücken, Germany

²Department of Computational Linguistics, Saarland University, Germany

³Department of Computer Science, Technische Universität München, Germany

Abstract. State-of-the-art human activity recognition methods build on discriminative learning which requires a representative training set for good performance. This leads to scalability issues for the recognition of large sets of highly diverse activities. In this paper we leverage the fact that many human activities are compositional and that the essential components of the activities can be obtained from textual descriptions or scripts. To share and transfer knowledge between composite activities we model them by a common set of attributes corresponding to basic actions and object participants. This attribute representation allows to incorporate script data that delivers new variations of a composite activity or even to unseen composite activities. In our experiments on 41 composite cooking tasks, we found that script data to successfully capture the high variability of composite activities. We show improvements in a supervised case where training data for all composite cooking tasks is available, but we are also able to recognize unseen composites by just using script data and without any manual video annotation.

1

Undoing the Damage of Dataset Bias

Aditya Khosla¹, Tinghui Zhou², Tomasz Malisiewicz¹,
Alexei A. Efros², and Antonio Torralba¹

¹Massachusetts Institute of Technology

{khosla,tomasz,torralba}@csail.mit.edu

²Carnegie Mellon University

{tinghuiz,efros}@cs.cmu.edu

Abstract. The presence of bias in existing object recognition datasets is now well-known in the computer vision community. While it remains in question whether creating an unbiased dataset is possible given limited resources, in this work we propose a discriminative framework that directly exploits dataset bias during training. In particular, our model learns two sets of weights: (1) bias vectors associated with each individual dataset, and (2) visual world weights that are common to all datasets, which are learned by undoing the associated bias from each dataset. The visual world weights are expected to be our best possible approximation to the object model trained on an unbiased dataset, and thus tend to have good generalization ability. We demonstrate the effectiveness of our model by applying the learned weights to a novel, unseen dataset, and report superior results for both classification and detection tasks compared to a classical SVM that does not account for the presence of bias. Overall, we find that it is beneficial to explicitly account for bias when combining multiple datasets.

1

Dog Breed Classification Using Part Localization

Jiongxin Li¹, Angjoo Kanazawa², David Jacobs², and Peter Belhumeur¹

¹Columbia University

²University of Maryland

Abstract. We propose a novel approach to fine-grained image classification in which instances from different classes share common parts but have wide variation in shape and appearance. We use dog breed identification as a test case to show that extracting corresponding parts improves classification performance. This domain is especially challenging since the appearance of corresponding parts can vary dramatically, e.g., the faces of bulldogs and beagles are very different. To find accurate correspondences, we build exemplar-based geometric and appearance models of dog breeds and their face parts. Part correspondence allows us to extract and compare descriptors in like image locations. Our approach also features a hierarchy of parts (e.g., face and eyes) and breed-specific part localization. We achieve 67% recognition rate on a large real-world dataset including 133 dog breeds and 8,351 images, and experimental results show that accurate part localization significantly increases classification performance compared to state-of-the-art approaches.

1

**A Dictionary Learning Approach for Classification:
Separating the Particularity and the Commonality**

Shu Kong and Donghui Wang

Dept. of Computer Science and Technology, Zhejiang University, Hangzhou, China
{aimerykong,dhwang}@zju.edu.cn

Abstract. Empirically, we find that, despite the class-specific features owned by the objects appearing in the images, the objects from different categories usually

share some common patterns, which do not contribute to the discrimination of them. Concentrating on this observation and under the general dictionary learning

(DL) framework, we propose a novel method to explicitly learn a common pattern pool (the commonality) and class-specific dictionaries (the particularity) for

classification. We call our method DL-COPAR, which can learn the most compact and most discriminative class-specific dictionaries used for classification. The proposed DL-COPAR is extensively evaluated both on synthetic data and on benchmark image databases in comparison with existing DL-based classification methods. The experimental results demonstrate that DL-COPAR achieves very promising performances in various applications, such as face recognition, handwritten digit recognition, scene classification and object recognition.

Keywords: Dictionary Learning, Classification, Commonality, Particularity.

1

Learning to Efficiently Detect Repeatable

Interest Points in Depth Data

Stefan Holzer^{1,2}, Jamie Shotton², and Pushmeet Kohli²

¹Department of Computer Science, CAMP, Technische Universität München (TUM)
holzers@in.tum.de

²Microsoft Research Cambridge

{Jamie.Shotton,pkohli}@microsoft.com

Abstract. Interest point (IP) detection is an important component of many computer vision methods. While there are a number of methods for detecting IPs in RGB images, modalities such as depth images and range scans have seen relatively little work. In this paper, we approach the IP detection problem from a machine learning viewpoint and formulate it as a regression problem. We learn a regression forest (RF) model that, given an image patch, tells us if there is an IP in the center of the patch. Our RF based method for IP detection allows an easy trade-off between speed and repeatability by adapting the depth and number of trees used for approximating the interest point response maps. The data used for training the RF model is obtained by running state-of-the-art IP detection methods

on the depth images. We show further how the IP response map used for training the RF can be specifically designed to increase repeatability by employing 3D models of scenes generated by reconstruction systems such as KinectFusion [1]. Our experiments demonstrate that the use of such data leads to considerably improved IP detection.

1

Effective Use of Frequent Itemset Mining for Image Classification

Basura Fernando¹, Elisa Fromont², and Tinne Tuytelaars¹

¹KU Leuven, ESAT-PSI, IBBT (Belgium)

²University of Saint-Etienne(France)

Abstract. In this paper we propose a new and effective scheme for applying frequent itemset mining to image classification tasks. We refer to the news extracted patterns as Frequent Local Histograms or FLHs.

During the construction of the FLHs, we pay special attention to keep all the local histogram information during the mining process and to select the most relevant reduced set of FLH patterns for classification. The careful choice of the visual primitives and some proposed extensions to exploit other visual cues such as colour or global spatial information allow us to build powerful bag-of-FLH-based image representations. We show that these bag-of-FLHs are more discriminative than traditional bag-of-words and yield state-of-the-art results on various image classification benchmarks.

1

Efficient Discriminative Projections for Compact Binary Descriptors

Tomasz Trzcinski and Vincent Lepetit

CVLab, EPFL, Lausanne, Switzerland

firstname.lastname@epfl.ch

Abstract. Binary descriptors of image patches are increasingly popular given that they require less storage and enable faster processing. This, however, comes at a price of lower recognition performances. To boost these performances, we project the image patches to a more discriminative subspace, and threshold their coordinates to build our binary descriptor. However, applying complex projections to the patches is slow, which negates some of the advantages of binary descriptors. Hence, our key idea is to learn the discriminative projections so that they can be decomposed into a small number of simple filters for which the responses can be computed fast. We show that with as few as 32 bits per descriptor we outperform the state-of-the-art binary descriptors in terms of both accuracy and efficiency.

1

Descriptor Learning Using Convex Optimisation

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman

Visual Geometry Group, University of Oxford

Abstract. The objective of this work is to learn descriptors suitable for the sparse feature detectors used in viewpoint invariant matching. We make a number of novel contributions towards this goal: First, it is shown that learning the pooling regions for the descriptor can be formulated as a convex optimisation problem selecting the regions using sparsity; second, it is shown that dimensionality reduction can also be formulated as a convex optimisation problem, using the nuclear norm to reduce dimensionality. Both of these problems use large margin discriminative learning methods. The third contribution is a new method of obtaining the positive and negative training data in a weakly supervised manner. And, finally, we employ a state-of-the-art stochastic optimizer that is efficient and well matched to the non-smooth cost functions proposed here. It is

demonstrated that the new learning methods improve over the state of the art in descriptor learning for large scale matching, Brown et al.[2], and large scale object retrieval, Philbin et al.[10].

1

Bottom-Up Perceptual Organization of Images
into Object Part Hypotheses

Maruthi Narayanan and Benjamin Kimia

Brown University

School of Engineering

Providence, RI 02912

{maruthi

narayanan,benjamin

kimia}@brown.edu

<http://vision.lems.brown.edu>

Abstract. The demise of "segmentation-then-recognition" strategy led to a paradigm shift toward feature-based discriminative recognition with significant success. However, increased complexity in multi-class datasets reveals that local low-level features may not be sufficiently discriminative, requiring the construction and use of more complex structural features which are necessarily category independent. The paper proposes a bottom-up procedure for generating fragment features which are intended to be object part hypotheses. Suggesting that the demise of segmentation to generate a representation suitable for recognition was due to prematurely committing to a grouping option in the face of ambiguities, the proposed framework considers and tracks multiple alternate grouping options. This approach is made tractable by (i)using amedial fragment representation which allows for the simultaneous use of multiple cues, (ii)a set of transforms to effect grouping operations, (iii) a containment graph representation which avoids duplicate consideration of possibilities, and the estimation of the likelihood of a grouping sequence to retain only plausible groupings. The resulting hypotheses are evaluated intrinsically by measuring their ability to represent objects with a few fragments. They are also evaluated by comparison to algorithms which aim to generate full object segments, with results that match or exceed the state of art, thus demonstrating the suitability of the proposed mid-level representation.

1

Match Graph Construction for Large Image Databases

K w a n g I n K i m l, James Tompkin^{1,2,3},

Martin Theobald¹, J a n K a u t z², and Christian Theobald¹

¹Max-Planck-Institut für Informatik, Campus E1 4, 66123 Saarbrücken, Germany

²University College London, Malet Place, WC1E 6BT London, UK

³Intel Visual Computing Institute, Campus E2 1, 66123 Saarbrücken, Germany

Abstract. How best to efficiently establish correspondence among a large set of images or video frames is an interesting unanswered question. For large databases, the high computational cost of performing pair-wise image matching is a major problem. However, for many applications, images are inherently sparsely connected, and so current techniques try to correctly estimate small potentially matching subsets of databases upon which to perform expensive pair-wise matching. Our contribution is to pose the identification of potential matches as a link

prediction problem in an image correspondence graph, and to propose an effective algorithm to solve this problem. Our algorithm facilitates incremental image

matching: initially, the match graph is very sparse, but it becomes dense as we al-

ternate between link prediction and verification. We demonstrate the effectiveness of our algorithm by comparing it with several existing alternatives on large-sc

ale

databases. Our resulting match graph is useful for many different applications. As an example, we show the benefits of our graph construction method to a label propagation application which propagates user-provided sparse object labels to other instances of that object in large image collections.

Keywords: Image matching, graph construction, link prediction.

1

Modeling Complex Temporal Composition of Actionlets for Activity Prediction

Kang Lil, Ji eHu2, a n dY u nF u l

1Department of ECE and College of CIS, Northeastern University, Boston, MA, USA

2Department of CSE, State University of New York, Buffalo, NY, USA

li.ka@husky.neu.edu, y.fu@neu.edu, jhu6@buffalo.edu

Abstract. Early prediction of ongoing activity has been more and more valuable in a large variety of time-critical applications. To build an effective representation for prediction, human activities can be characterized by a complex temporal composition of constituent simple actions. Different from early recognition on short-duration simple activities, we propose a novel framework for long-duration complex activity prediction

by discovering the causal relationships between constituent actions and the predictable characteristics of activities. The major contributions of our work include: (1) we propose a novel activity decomposition method by monitoring motion velocity which encodes a temporal decomposition of long activities into a sequence of meaningful action units; (2) Probabilistic Suffix Tree (PST) is introduced to represent both large and small order Markov dependencies between action units; (3) we present a Predictive Accumulative Function (PAF) to depict the predictability of each kind of activity. The effectiveness of the proposed method is evaluated on two experimental scenarios: activities with middle-level complexity and activities with high-level complexity. Our method achieves promising results and can predict global activity classes and local action units.

1

Learning Human Interaction

by Interactive Phrases

Yu Kong1,3, Yunde Jia1, a n dY u nF u2

1Beijing Laboratory of Intelligent Information Technology

School of Computer Science, Beijing Institute of Technology

Beijing 100081, P.R. China

2Department of ECE and College of CIS, Northeastern University, Boston, MA

3Department of CSE, State University of New York, Buffalo, NY

{kongyu, jiaiyunde}@bit.edu.cn, y.fu@neu.edu

Abstract. In this paper, we present a novel approach for human interaction recognition from videos. We introduce high-level descriptions called interactive phrases to express binary semantic motion relationships between interacting people. Interactive phrases naturally exploit human knowledge to describe interactions and allow us to construct a more descriptive model for recognizing human interactions. We propose a novel hierarchical model to encode interactive phrases based on the latent SVM framework where interactive phrases are treated as latent variables. The interdependencies between interactive phrases are explicitly captured in the model to deal with motion ambiguity and partial occlusion in interactions. We evaluate our method on a newly collected BIT-Interaction dataset and UT-Interaction dataset. Promising results demonstrate the effectiveness of the proposed method.

1

Learning to Recognize Daily Actions Using Gaze

Alireza Fathi, Yin Li, and James M. Rehg

College of Computing
Georgia Institute of Technology

Abstract. We present a probabilistic generative model for simultaneously recognizing daily actions and predicting gaze locations in videos recorded from an egocentric camera. We focus on activities requiring eye-hand coordination and model the spatio-temporal relationship between the gaze point, the scene objects, and the action label. Our model captures the fact that the distribution of both visual features and object occurrences in the vicinity of the gaze point is correlated with the verb-object pair describing the action. It explicitly incorporates known properties of gaze behavior from the psychology literature, such as the temporal delay between fixation and manipulation events. We present an inference method that can predict the best sequence of gaze locations and the associated action label from an input sequence of images. We demonstrate improvements in action recognition rates and gaze prediction accuracy relative to state-of-the-art methods, on two new datasets that contain egocentric videos of daily activities and gaze.

1

Gait Recognition by Ranking

Ra'ul Mart'ın-F'elez¹ and Tao Xiang²

¹Institute of New Imaging Technologies, Universitat Jaume I, Castell' o 12071, Spain

²School of EECS, Queen Mary, University of London, London E1 4NS, U.K.

martinr@uji.es, txiang@eecs.qmul.ac.uk

Abstract. The advantage of gait over other biometrics such as face or fingerprint is that it can operate from a distance and without subject cooperation. However, this also makes gait subject to changes in various covariate conditions including carrying, clothing, surface and view angle. Existing approaches attempt to address these condition changes by feature selection, feature transformation or discriminant subspace learning. However, they suffer from lack of training samples from each subject, can only cope with changes in a subset of conditions with limited success, and are based on the invalid assumption that the covariate conditions are known a priori. They are thus unable to perform gait recognition under a genuine uncooperative setting. We propose a novel approach which casts gait recognition as a bipartite ranking problem and leverages training samples from different classes/people and even from different datasets. This makes our approach suitable for recognition under a genuine uncooperative setting and robust against any covariate types, as demonstrated by our extensive experiments.

Keywords: Gait recognition, Learning to rank, Transfer learning.

1

Semi-intrinsic Mean Shift on Riemannian Manifolds

Rui Caseiro, Jo'ao F. Henriques, Pedro Martins, and Jorge Batista

Institute of Systems and Robotics - University of Coimbra, Portugal

[ruicaseiro](mailto:ruicaseiro@isr.uc.pt), [henriques](mailto:henriques@isr.uc.pt), [pedromartins](mailto:pedromartins@isr.uc.pt), [batista](mailto:batista@isr.uc.pt)

Abstract. The original mean shift algorithm [1] on Euclidean spaces (MS) was extended in [2] to operate on general Riemannian manifolds. This extension is extrinsic (Ext-MS) since the mode seeking is performed on the tangent spaces [3], where the underlying curvature is not fully considered (tangent spaces are only valid in a small neighborhood). In [3] was proposed an intrinsic mean shift designed to operate on two particular Riemannian manifolds (IntGS-MS), i.e. Grassmann and Stiefel manifolds (using manifold-dedicated density kernels). It is then natural to ask whether mean shift could be intrinsically extended to work on a large class of manifolds. We propose a novel paradigm to intrinsically reformulate the mean shift on general Riemannian manifolds. This is ac-

complished by embedding the Riemannian manifold into a Reproducing Kernel Hilbert Space (RKHS) by using a general and mathematically well-founded Riemannian kernel function, i.e. heat kernel [4]. The key issue is that when the data is implicitly mapped to the Hilbert space, the curvature of the manifold is taken into account (i.e. exploits the underlying information of the data). The inherent optimization is then performed on the embedded space. Theoretic analysis and experimental results demonstrate the promise and effectiveness of this novel paradigm.

1

Efficient Nonlocal Regularization

for Optical Flow

Philipp Krähenbühl and Vladlen Koltun

Stanford University

{philkr,vladlen}@cs.stanford.edu

Abstract. Dense optical flow estimation in images is a challenging problem because the algorithm must coordinate the estimated motion across large regions in the image, while avoiding inappropriate smoothing over motion boundaries. Recent works have advocated for the use of nonlocal regularization to model long-range correlations in the flow. However, incorporating nonlocal regularization into an energy optimization framework is challenging due to the large number of pairwise penalty terms. Existing techniques either substitute intermediate filtering of the flow field for direct optimization of the nonlocal objective, or suffer substantial performance penalties when the range of the regularizer increases. In this paper, we describe an optimization algorithm that efficiently handles a general type of nonlocal regularization objectives for optical flow estimation. The computational complexity of the algorithm is independent of the range of the regularizer. We show that nonlocal regularization improves estimation accuracy at longer ranges than previously reported, and is complementary to intermediate filtering of the flow field. Our algorithm is simple and is compatible with many optical flow models.

1

Fast Fusion Moves for Multi-model Estimation

Andrew DeLong, Olga Veksler, and Yuri Boykov

University of Western Ontario, Canada

Abstract. We develop a fast, effective algorithm for minimizing a well-known objective function for robust multi-model estimation. Our work introduces a com-

binatorial step belonging to a family of powerful move-making methods like α -expansion and fusion. We also show that our subproblem can be quickly transformed into a comparatively small instance of minimum-weighted vertex-cover. In practice, these vertex-cover subproblems are almost always bipartite and can be solved exactly by specialized network flow algorithms. Experiments indicate that our approach achieves the robustness of methods like affinity propagation, whilst providing the speed of fast greedy heuristics.

1

Approximate MRF Inference Using

Bounded Treewidth Subgraphs

Alexander Fixl, Joyce Chen¹, Endre Boros², and Ramnath A. Iyer

¹Cornell University, Computer Science Department, Ithaca, New York

{afix,yuhsin,rdz}@cs.cornell.edu

²Rutgers University, RUTCOR, New Brunswick, New Jersey

boros@rci.rutgers.edu

Abstract. Graph cut algorithms [9], commonly used in computer vision, solve a first-order MRF over binary variables. The state of the art for this NP-hard problem is QPBO [1,2], which finds the values for a subset of the variables in the

global minimum. While QPBO is very effective overall there are still many difficult problems where it can only label a small subset of the variables. We propose a new approach that, instead of optimizing the original graphical model, instead

optimizes a tractable sub-model, defined as an energy function that uses a subset

of the pairwise interactions of the original, but for which exact inference can be done efficiently. Our Bounded Treewidth Subgraph (k-BTS) algorithm greedily computes a large weight treewidth-k subgraph of the signed graph, then solves the energy minimization problem for this subgraph by dynamic programming. The edges omitted by our greedy method provide a per-instance lower bound. We demonstrate promising experimental results for binary deconvolution, a challenging problem used to benchmark QPBO [2]: our algorithm performs an order of magnitude better than QPBO or its common variants [4], both in terms of energy and accuracy, and the visual quality of our output is strikingly better as well.

We also obtain a significant improvement in energy and accuracy on a stereo benchmark with 2nd order priors [5], although the improvement in visual quality is more modest. Our method's running time is comparable to QPBO.

1

Recursive Bilateral Filtering

Qingxiong Yang

Department of Computer Science,

City University of Hong Kong, Hong Kong, China

<http://www.cs.cityu.edu.hk/~qiyang/>

Abstract. This paper proposes a recursive implementation of the bilateral filter. Unlike previous methods, this implementation yields an bilateral filter whose computational complexity is linear in both input size and dimensionality. The proposed implementation demonstrates that the bilateral filter can be as efficient as the recent edge-preserving filtering methods, especially for high-dimensional images. Let the number of pixels contained in the image be N , and the number of channels be D , the computational complexity of the proposed implementation will be $O(ND)$. It is more efficient than the state-of-the-art bilateral filtering methods that have a computational complexity of $O(ND^2)$ [1] (linear in the image size but polynomial in dimensionality) or $O(N \log(N)D)$ [2] (linear in the dimensionality thus faster than [1] for high-dimensional filtering). Specifically, the proposed implementation takes about 43 ms to process a one megapixel color image (and about 14 ms to process a 1 megapixel grayscale image) which is about 18× faster than [1] and 86× faster than [2]. The experiments were conducted on a MacBook Air laptop computer with a 1.8 GHz Intel Core i7 CPU and 4 GB memory. The memory complexity of the proposed implementation is also low: as few as the image memory will be required (memory for the images before and after filtering is excluded). This paper also derives a new filter named gradient domain bilateral filter from the proposed recursive implementation. Unlike the bilateral filter, it performs bilateral filtering on the gradient domain. It can be used for edge-preserving filtering but avoids sharp edges that are observed to cause visible artifacts in some computer graphics tasks. The proposed implementations were proved to be effective for a number of computer vision and computer graphics applications, including stylization, tone mapping, detail enhancement and stereo matching.

1

Accelerated Large Scale Optimization

by Concomitant Hashing

Yadong Mu, John Wright, and Shih-Fu Chang

Electrical Engineering Department,

Columbia University, New York, NY 10027

{muyadong,johnwright,sfchang}@ee.columbia.edu

Abstract. Traditional locality-sensitive hashing (LSH) techniques aim to tackle the curse of explosive data scale by guaranteeing that similar samples are projected onto proximal hash buckets. Despite the success of LSH on numerous vision tasks like image retrieval and object matching, however, its potential in large-scale optimization is only realized recently. In this paper we further advance this nascent area. We first identify two common operations known as the computational bottleneck of numerous optimization algorithms in a large-scale setting, i.e., min/max inner product. We propose a hashing scheme for accelerating min/max inner product, which exploits properties of order statistics of statistically correlated random vectors. Compared with other schemes, our algorithm exhibits improved recall at a lower computational cost. The effectiveness and efficiency of the proposed method are corroborated by theoretic analysis and several important applications. Especially, we use the proposed hashing scheme to perform approximate ℓ_2 -regularized least squares with dictionaries with millions of elements, a scale which is beyond the capability of currently known exact solvers. Nonetheless, it is highlighted that the focus of this paper is not on a new hashing scheme for approximate nearest neighbor problem. It exploits a new application for the hashing techniques and proposes a general framework for accelerating a large variety of optimization procedures in computer vision.

1

Graph Degree Linkage:

Agglomerative Clustering on a Directed Graph

Wei Zhang¹, Xiaogang Wang^{2,3}, Delia Zha¹, and Xiaoou Tang^{1,3}

¹Department of Information Engineering, The Chinese University of Hong Kong
wzhang009@gmail.com

²Department of Electronic Engineering, The Chinese University of Hong Kong

³Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, China

Abstract. This paper proposes a simple but effective graph-based agglomerative algorithm, for clustering high-dimensional data. We explore the different roles of

two fundamental concepts in graph theory, indegree and outdegree, in the context of clustering. The average indegree reflects the density near a sample, and the average outdegree characterizes the local geometry around a sample. Based on such insights, we define the affinity measure of clusters via the product of average indegree and average outdegree. The product-based affinity makes our algorithm robust to noise. The algorithm has three main advantages: good performance, easy implementation, and high computational efficiency. We test the algorithm on two fundamental computer vision problems: image clustering and object matching. Extensive experiments demonstrate that it outperforms the state-

of-the-arts in both applications.

1

Supervised Earth Mover's Distance Learning

and Its Computer Vision Applications

Fan Wang and Leonidas J. Guibas

Stanford University, CA, United States

Abstract. The Earth Mover's Distance (EMD) is an intuitive and natural distance metric for comparing two histograms or probability distributions.

It provides a distance value as well as a flow-network indicating how the probability mass is optimally transported between the bins. In traditional EMD, the ground distance between the bins is pre-defined. Instead, we propose to jointly optimize the ground distance matrix and the EMD flow-network based on a partial ordering of histogram dis-

tances in an optimization framework. Our method is further extended to accept information from general labeled pairs. The trained ground distance better reflects the cross-bin relationships, hence produces more accurate EMD values and flow-networks. Two computer vision applications are used to demonstrate the effectiveness of the algorithm: first, we apply the optimized EMD value to face verification, and achieve state-of-the-art performance on the PubFig and the LFW data sets; second, the learned EMD flow-network is used to analyze face attribute changes, obtaining consistent paths that demonstrate intuitive transitions on certain facial attributes.

1

Global Optimization of Object Pose and Motion
from a Single Rolling Shutter Image
with Automatic 2D-3D Matching

Ludovic Magerand^{1,2}, Adrien Bartoli², Omar Ait-Aider¹, and Daniel Pizarro³

¹Institut Pascal – Université Blaise Pascal – Clermont-Ferrand

²ISIT – Université d’Auvergne – Clermont-Ferrand

³University of Alcala – Alcala de Henares

Abstract. Low cost CMOS cameras can have an acquisition mode called rolling shutter which sequentially exposes the scan-lines. When a single object moves with respect to the camera, this creates image distortions. Assuming 2D-3D correspondences known, previous work showed that the object pose and kinematics can be estimated from a single rolling shutter image. This was achieved using a suboptimal initialization followed by local iterative optimization.

We propose a polynomial projection model for rolling shutter cameras and a constrained global optimization of its parameters. This is done by means of a semidefinite programming problem obtained from the generalized problem of moments method. Contrarily to previous work, our optimization does not require an initialization and ensures that the global minimum is achieved. This allows us to build automatically robust 2D-3D correspondences using a template to provide an initial set of correspondences.

Experiments show that our method slightly improves previous work on both simulated and real data. This is due to local minima into which previous methods get trapped. We also successfully experimented building 2D-3D correspondences automatically with both simulated and real data.

Keywords: rolling shutter, motion estimation, matching.

1

Online Learning of Linear Predictors
for Real-Time Tracking

Stefan Holzer¹, Marc Pollefeys²,

Slobodan Ilic¹, David Joseph Tan¹, and Nassir Navab¹

¹Department of Computer Science, Technische Universität München (TUM), Boltzmannstrasse 3, 85748 Garching, Germany

{holzers,slobodan.ilic,tanda,navab}@in.tum.de

²Department of Computer Science, ETH Zurich, CNB G105,

Universitätsstrasse 6, CH-8092 Zurich, Switzerland

marc.pollefeys@inf.ethz.ch

Abstract. Although fast and reliable, real-time template tracking using linear predictors requires a long training time. The lack of the ability to learn new templates online prevents their use in applications that require fast learning. This especially holds for applications where the scene is not known a priori and multiple templates have to be added online. So far, linear predictors had to be either learned offline [1] or in an iterative manner by starting with a small sized template and growing it overtime [2]. In this paper, we propose a fast and simple reformulation of the learning procedure that allows learning new linear predictors online.

Keywords: template tracking, template learning, linear predictors.

1

Online Learned Discriminative Part-Based Appearance Models for Multi-human Tracking

Bo Yang and Ram Nevatia

Institute for Robotics and Intelligent Systems,
University of Southern California

Los Angeles, CA 90089, USA

{yangbo,nevatia}@usc.edu

Abstract. We introduce an online learning approach to produce discriminative part-based appearance models (DPAMs) for tracking multiple humans in real scenes by incorporating association based and category free tracking methods. Detection responses are gradually associated into tracklets in multiple levels to produce final tracks. Unlike most previous multi-target tracking approaches which do not explicitly consider occlusions in appearance modeling, we introduce a part based model that explicitly finds unoccluded parts by occlusion reasoning in each frame, so that occluded parts are removed in appearance modeling. Then DPAMs for each tracklet is online learned to distinguish a tracklet with others as well as the background, and is further used in a conservative category free tracking approach to partially overcome the missed detection problem as well as to reduce difficulties in tracklet associations under long gaps. We evaluate our approach on three public data sets, and show significant improvements compared with state-of-art methods.

Keywords: multi-human tracking, online learned discriminative models.

1

Exposure Stacks of Live Scenes with Hand-Held Cameras

Jun Hu¹, Orazio Gallo², and Kari Pulli²

¹Department of Computer Science, Duke University

²NVIDIA Research, Santa Clara, CA

Abstract. Many computational photography applications require the user to take multiple pictures of the same scene with different camera settings. While this allows to capture more information about the scene than what is possible with a single image, the approach is limited by the requirement that the images be perfectly registered. In a typical scenario the camera is hand-held and is therefore prone to moving during the capture of an image burst, while the scene is likely to contain moving objects. Combining such images without careful registration introduces annoying artifacts in the final image. This paper presents a method to register exposure stacks in the presence of both camera motion and scene changes.

Our approach warps and modifies the content of the images in the stack to match that of a reference image. Even in the presence of large, highly non-rigid displacements we show that the images are correctly registered to the reference.

1

Dual-Force Metric Learning for Robust Distracter-Resistant Tracker

Zhibin Hong¹, Xueming Li², and Dacheng Tao¹

¹Centre for Quantum Computation and Intelligent Systems,

Faculty of Engineering and Information Technology,

University of Technology, Sydney, NSW, Australia

²Center for Automation Research, Electrical & Computer Engineering Department,

University of Maryland, College Park, MD, USA

Abstract. In this paper, we propose a robust distracter-resistant tracking approach by learning a discriminative metric that adaptively learns the importance of features on-the-fly. The proposed metric is elaborately designed

for the tracking problem by forming a margin objective function which systematically includes distance margin maximization and reconstruction error constraint that acts as a force to push distracters away from the positive space and into the negative space. Due to the variety of negative samples in the tracking problem, we specifically introduce the similarity propagation technique that gives distracters a second force from the negative space. Consequently, the discriminative metric obtained helps to preserve the most discriminative information to separate the target from distracters while ensuring the stability of the optimal metric. We seamlessly combine it with the popular L1 minimization tracker. Our tracker is therefore not only resistant to distracters, but also inherits the merit of occlusion robustness from the L1 tracker. Quantitative comparisons with several state-of-the-art algorithms have been conducted in many challenging video sequences. The results show that our method resists distracters excellently and achieves superior performance.

Keywords: Visual tracking, distracter, distance metric, similarity propagation.

1

Shape and Reflectance from Natural Illumination

Geoffrey Oxholm and Ko Nishino

Department of Computer Science

Drexel University, Philadelphia, PA 19104, USA

{gao25,kon}@drexel.edu

Abstract. We introduce a method to jointly estimate the BRDF and geometry of an object from a single image under known, but uncontrolled, natural illumination. We show that this previously unexplored problem becomes tractable when one exploits the orientation clues embedded in the lighting environment. Intuitively, unique regions in the lighting environment act analogously to the point light sources of traditional photometric stereo; they strongly constrain the orientation of the surface patches that reflect them. The reflectance, which acts as a bandpass filter on the lighting environment, determines the necessary scale of such regions. Accurate reflectance estimation, however, relies on accurate surface orientation information. Thus, these two factors must be estimated jointly. To do so, we derive a probabilistic formulation and introduce priors to address situations where the reflectance and lighting environment do not sufficiently constrain the geometry of the object. Through extensive experimentation we show what this space looks like, and offer insights into what problems become solvable in various categories of real-world natural illumination environments.

1

entation of the surface patches that reflect them. The reflectance, which acts as a bandpass filter on the lighting environment, determines the necessary scale of such regions. Accurate reflectance estimation, however, relies on accurate surface orientation information. Thus, these two factors must be estimated jointly. To do so, we derive a probabilistic formulation and introduce priors to address situations where the reflectance and lighting environment do not sufficiently constrain the geometry of the object. Through extensive experimentation we show what this space looks like, and offer insights into what problems become solvable in various categories of real-world natural illumination environments.

1

so, we derive a probabilistic formulation and introduce priors to address situations where the reflectance and lighting environment do not sufficiently constrain the geometry of the object. Through extensive experimentation we show what this space looks like, and offer insights into what problems become solvable in various categories of real-world natural illumination environments.

1

Frequency Analysis of Transient Light Transport

with Applications in Bare Sensor Imaging

Di Wu^{1,2,5,*}, Gordon Wetzstein¹, Christopher Barsil¹, Thomas Willwacher³,

Matthew O'Toole⁴, Nikhil Nalwa¹, Qionghai Dai²,

Kyros Kutulakos⁴, and Ramesh Raskar¹

¹MIT Media Lab

²Department of Automation, Tsinghua University

³Department of Mathematics, Harvard University

⁴Department of Computer Science, University of Toronto

⁵Graduate School at Shenzhen, Tsinghua University

Abstract. Light transport has been analyzed extensively, in both the primal domain and the frequency domain; the latter provides intuition of effects introduced by free space propagation and by optical elements, and allows for optimal designs of computational cameras for tailored, efficient information capture. Here, we relax the common assumption that the speed of light is infinite and analyze free space propagation in

the frequency domain considering spatial, temporal, and angular light variation. Using this analysis, we derive analytic expressions for cross-dimensional information transfer and show how this can be exploited for designing a new, time-resolved bare sensor imaging system.

Keywords: Light transport, Fourier analysis, Time of Flight, Lensless imaging.

1

Nonuniform Lattice Regression

for Modeling the Camera Imaging Pipeline

Hai Ting Lin¹, Zheng Lu²,

Seon Joo Kim³, and Michael S. Brown¹

¹National University of Singapore

²University of Texas at Austin

³SUNY Korea

Abstract. We describe a method to construct a lookup table

(LUT) that is effective in modeling the camera imaging pipeline that maps a RAW camera values to their sRGB output. This work builds on the recent in-camera color processing model proposed by Kim et al. [1] that included a 3D gamut-mapping function. The major drawback in [1] is the high computational cost of the 3D mapping function that uses radial basis functions (RBF) involving several thousand control points. We show how to construct a LUT using a novel nonuniform lattice regression method that adapts the LUT lattice to better fit the 3D gamut-mapping function. Our method offers not only a performance speedup of an order of magnitude faster than RBF, but also a compact mechanism to describe the imaging pipeline.

1

Context-Based Automatic Local

Image Enhancement

Sung Ju Hwang¹, Ashish Kapoor², and Sing Bing Kang²

¹The University of Texas, Austin, TX, USA

sjhwang@cs.utexas.edu

²Microsoft Research, Redmond, WA, USA

{akapoor, sbkang}@microsoft.com

Abstract. In this paper, we describe a technique to automatically enhance the perceptual quality of an image. Unlike previous techniques, where global statistics of the image are used to determine enhancement operation, our method is local and relies on local scene descriptors and context in addition to high-level image statistics. We cast the problem of image enhancement as searching for the best transformation for each pixel in the given image and then discovering the enhanced image using a formulation based on Gaussian Random Fields. The search is done in a coarse-to-fine manner, namely by finding the best candidate images, followed by pixels. Our experiments indicate that such context-based local enhancement is better than global enhancement schemes. A user study using Mechanical Turk shows that the subjects prefer contextual and local enhancements over the ones provided by existing schemes.

1

Segmentation with Non-linear Regional

Constraints via Line-Search Cuts

Lena Gorelick¹, Frank R. Schmidt², Yuri Boykov¹,

Andrew Delong¹, and Aaron Ward¹

¹University of Western Ontario, Canada

²Universit e Paris Est, France

Abstract. This paper is concerned with energy-based image segmentation problems. We introduce a general class of regional functionals

defined as an arbitrary non-linear combination of regional unary terms. Such (high-order) functionals are very useful in vision and medical applications and some special cases appear in prior art. For example, our general class of functionals includes but is not restricted to soft constraints on segment volume, its appearance histogram, or shape. Our overall segmentation energy combines regional functionals with standard length-based regularizers and/or other submodular terms. In general, regional functionals make the corresponding energy minimization NP-hard. We propose a new greedy algorithm based on iterative line search. A parametric max-flow technique efficiently explores all solutions along the direction (line) of the steepest descent of the energy. We compute the best "step size", i.e. the globally optimal solution along the line. This algorithm can make large moves escaping weak local minima, as demonstrated on many real images.

1

Hausdorff Distance Constraint
for Multi-surface Segmentation
Frank R. Schmidtland Yuri Boykov²
¹Universit e Paris Est, France
²University of Western Ontario, Canada

Abstract. It is well known that multi-surface segmentation can be cast as a multi-labeling problem. Different segments may belong to the same semantic object which may impose various inter-segment constraints [1]. In medical applications, there are a lot of scenarios where upper bounds on the Hausdorff distances between subsequent surfaces are known. We show that incorporating these priors into multi-surface segmentation is potentially NP-hard. To cope with this problem we develop a submodular-supermodular procedure that converges to a locally optimal solution well-approximating the problem. While we cannot guarantee global optimality, only feasible solutions are considered during the optimization process. Empirically, we get useful solutions for many challenging medical applications including MRI and ultrasound images.

1

Background Subtraction Using Low Rank
and Group Sparsity Constraints
Xinyi Cui¹, Junzhou Huang², Shaoting Zhang¹, and Dimitris N. Metaxas¹
¹CS Dept., Rutgers University
Piscataway, NJ 08854, USA
²CSE Dept., Univ. of Texas at Arlington
Arlington, TX, 76019, USA
{xycui, shaoting, dnm}@cs.rutgers.edu,
jzhuang@uta.edu

Abstract. Background subtraction has been widely investigated in recent years. Most previous work has focused on stationary cameras. Recently, moving cameras have also been studied since videos from mobile devices have increased significantly. In this paper, we propose a unified and robust framework to effectively handle diverse types of videos, e.g., videos from stationary or moving cameras. Our model is inspired by two observations: 1) background motion caused by orthographic cameras lies in a low rank subspace, and 2) pixels belonging to one trajectory tend to group together. Based on these two observations, we introduce a new model using both low rank and group sparsity constraints. It is able to robustly decompose a motion trajectory matrix into foreground and background ones. After obtaining foreground and background trajectories, the information gathered on them is used to build a statistical model to further label frames at the pixel level. Extensive experiments demonstrate very competitive performance on both synthetic data and real videos.

1

FreeHand - Drawn Sketch Segmentation

Zhenbang Sun¹, Changhu Wang², Liqing Zhang¹, and Lei Zhang²

¹Brain-Like Computing Lab, Shanghai Jiao Tong University, P.R. China

²Microsoft Research Asia

Abstract. In this paper, we study the problem of how to segment a freehand sketch at the object level. By carefully considering the basic principles of human perceptual organization, a real-time solution is presented to automatically segment a user's sketch during his/her drawing. First, a graph-based sketch segmentation algorithm is proposed to segment a cluttered sketch into multiple parts based on the factor of proximity. Then, to improve the ability of detecting semantically meaningful objects, a semantic-based approach is introduced to simulate the past experience in the perceptual system by leveraging a web-scale clipart database. Finally, other important factors learnt from past experience, such as similarity, symmetry, direction, and closure, are taken into account to make the approach more robust and practical. The proposed sketch segmentation framework has ability to handle complex sketches with overlapped objects. Extensive experimental results show the effectiveness of the proposed framework and algorithms.

1

Auto-Grouped Sparse Representation

for Visual Analysis

Jiashi Feng¹, Xiaotong Yuan², Zilei Wang³, Hua Xu⁴, and Shuicheng Yan¹

¹Department of ECE, National University of Singapore

²Department of Statistics, Rutgers University

³Department of Automation, University of Science and Technology of China

⁴Department of ME, National University of Singapore

Abstract. In this work, we investigate how to automatically uncover the underlying group structure of a feature vector such that each group characterizes certain object-specific patterns, e.g., visual patterns. In this paper, we propose a novel auto-grouped sparse representation (ASR) method. ASR groups semantically correlated feature elements together through optimally fusing their multiple sparse representations. Due to the intractability of primal objective function, we also propose well-behaved convex relaxation and smooth approximation to guarantee obtaining a global optimal solution effectively. Finally, we apply ASR in two important visual analysis tasks: multi-label image classification and motion segmentation. Comprehensive experimental evaluations show that ASR is able to achieve superior performance compared with the state-of-the-arts on these two tasks.

1

A QCQP Approach to Triangulation

Chris Aholt¹, Sameer Agarwal², and Rekha Thomas¹

¹University of Washington

²Google Inc.

Abstract. Triangulation of a three-dimensional point from two-dimensional images can be formulated as a quadratically constrained quadratic program. We propose an algorithm to extract candidate solutions to this problem from its semidefinite programming relaxations. We then describe a sufficient condition and a polynomial time test for certifying when such a solution is optimal. This test has no false positives.

Experiments indicate that false negatives are rare, and the algorithm has excellent performance in practice. We explain this phenomenon in terms of the geometry of the triangulation problem.

1

Reconstructing the World's Museums

Jianxiong Xiao¹ and Yasutaka Furukawa²

¹Massachusetts Institute of Technology

²Google Inc.

Abstract. Photorealistic maps are a useful navigational guide for large indoor environments, such as museums and businesses. However, it is impossible to acquire photographs covering a large indoor environment from aerial viewpoints. This paper presents a 3D reconstruction and visualization system to automatically produce clean and well-regularized texture-mapped 3D models for large indoor scenes, from ground-level photographs and 3D laser points. The key component is a new algorithm called "Inverse CSG" for reconstructing a scene in a Constructive Solid Geometry (CSG) representation consisting of volumetric primitives, which imposes powerful regularization constraints to exploit structural regularities. We also propose several techniques to adjust the 3D model to make it suitable for rendering the 3D maps from aerial viewpoints. The visualization system enables users to easily browse a large scale indoor environment from a bird's-eye view, locate specific room interiors, fly into a place of interest, view immersive ground-

level panorama views, and zoom out again, all with seamless 3D transitions. We demonstrate our system on various museums, including the Metropolitan Museum of Art in New York City - one of the largest art galleries in the world.

1

Background Inpainting for Videos with Dynamic

Objects and a Free-Moving Camera

Miguel Granados¹, Kwangmin Kim¹, James Tompkins^{1,2,3},

Jan Kautz², and Christian Theobalt¹

¹Max-Planck-Institut für Informatik, Campus E1 4, 66123 Saarbrücken, Germany

²University College London, Malet Place, WC1E 6BT London, UK

³Intel Visual Computing Institute, Campus E2 1, 66123 Saarbrücken, Germany

Abstract. We propose a method for removing marked dynamic objects from videos captured with a free-moving camera, so long as the objects occlude parts of the scene with a static background. Our approach takes as input a video, a mask marking the object to be removed, and a mask marking the dynamic objects to remain in the scene. To inpaint a frame, we align other candidate frames in which parts of the missing region are visible.

Among these candidates, a single source is chosen to fill each pixel so that the final arrangement is color-consistent. Intensity differences between sources are smoothed using gradient domain fusion. Our frame alignment process assumes that the scene can be approximated using piecewise planar geometry: A set of homographies is estimated for each frame pair, and one each is selected for aligning pixels such that the color-discrepancy is minimized and the epipolar constraints are maintained. We provide experimental validation with several real-world video sequences to demonstrate that, unlike in previous work, inpainting video shot with free-moving cameras does not necessarily require estimation of absolute camera positions and per-frame per-pixel depth maps.

Keywords: video processing, video completion, video inpainting, image alignment, background estimation, free-camera, graph-cuts.

1

Optimal Templates for Nonrigid Surface

Reconstruction

Markus Moll¹ and Luc Van Gool^{1,2}

¹ESAT-IBBT/PSI - KU Leuven

{markus.moll,luc.vangool }@esat.kuleuven.be

2Computer Vision Laboratory - Swiss Federal Institute of Technology

vangool@vision.ee.ethz.ch

Abstract. This paper addresses the problem of reconstructing a de-
forming surface from point observations in a monocular video sequence. Recent state-of-the-art approaches divide the surface into smaller patches to simplify the problem. Among these, one very promising approach reconstructs the patches individually using a quadratic deformation model. In this paper, we demonstrate limitations that affect its applicability to real-world data and propose an approach that overcomes these problems. In particular, we show how to eliminate the need for manually picking a template that is used to model the deformations. We evaluate our algorithm on both synthetic and real-world data sets and show that it systematically reduces the reconstruction error by a factor of up to ten.

Fig. 1. A textured waving flag with overlaid vertex mesh of which both the shape and

the deformations are recovered

1

Learning Domain Knowledge

for Façade Labelling

Dengxin Dai^{1,2}, Mukta Prasad¹, Gerhard Schmitt², and Luc Van Gool¹

¹Computer Vision Lab, ETH Zürich

²Chair for Information Architecture, ETH Zürich

Abstract. This paper presents an approach to address the problem of image façade labelling. In the architectural literature, domain knowledge is usually expressed geometrically in the formal design, so façade labelling should on the one hand conform to visual evidence, and on the other hand to the architectural principles – how individual assets (e.g. doors, windows) interact with each other to form a façade as a whole. To this end, we first propose a recursive splitting method to segment façades into a bunch of tiles for semantic recognition. The segmentation improves the processing speed, guides visual recognition on suitable scales and renders the extraction of architectural principles easy. Given a set of segmented training façades with their label maps, we then identify a set of meta-features to capture both the visual evidence and the architectural principles. The features are used to train our façade labelling model. In the test stage, the features are extracted from segmented façades and the inferred label maps. The following three steps are iterated until the optimal labelling is reached: 1) proposing modifications to the current labelling; 2) extracting new features for the proposed labelling; 3) feeding the new features to the labelling model to decide whether to accept the modifications. In experiments, we evaluated our method on the ECP façade dataset and achieved higher precision than the state-of-the-art at both the pixel level and the structural level.

1

Simultaneous Shape and Pose Adaption

of Articulated Models Using Linear

Optimization

Matthias Straka, Stefan Hauswiesner, Matthias Ruether, and Horst Bischof
Institute for Computer Graphics and Vision, Graz University of Technology,
Inffeldgasse 16/II, A-8010 Graz, Austria

{straka,hauswiesner,ruether,bischof }@icg.tugraz.at

<http://www.icg.tugraz.at/>

Abstract. We propose a novel formulation to express the attachment of a polygonal surface to a skeleton using purely linear terms. This enables to simultaneously adapt the pose and shape of an articulated model in an efficient way. Our work is motivated by the difficulty to constrain a mesh when adapting it to multi-view silhouette images. However, such

an adaption is essential when capturing the detailed temporal evolution of skin and clothing of a human actor without markers. While related work is only able to ensure surface consistency during mesh adaption, our coupled optimization of the skeleton creates structural stability and minimizes the sensibility to occlusions and outliers in input images. We demonstrate the benefits of our approach in an extensive evaluation. The skeleton attachment considerably reduces implausible deformations, especially when the number of input views is limited.

Keywords: Shape Adaption, Pose Estimation, Mesh Editing, Linear Optimization.

1

Robust Fitting for Multiple View Geometry

Olof Engvist, Erik Ask, Fredrik Kahl, and Kalle Åström

Centre for Mathematical Sciences, Lund University

<http://www.maths.lth.se/vision/>

Abstract. How hard are geometric vision problems with outliers? We show that for most fitting problems, a solution that minimizes the number of outliers can be found with an algorithm that has polynomial time-complexity in the number of points (independent of the rate of outliers). Further, and perhaps more interestingly, other cost functions such as the truncated L

2-norm can also be handled within the same framework with the same time complexity.

We apply our framework to triangulation, relative pose problems and stitching, and give several other examples that fulfill the required conditions.

Based on efficient polynomial equation solvers, it is experimentally demonstrated that these problems can be solved reliably, in particular for low-dimensional models. Comparisons to standard random sampling solvers are also given.

1

Improving Image-Based Localization by Active

Correspondence Search

Torsten Sattler¹, Bastian Leibe², and Leif Kobbelt¹

¹RWTH Aachen University, Aachen, Germany

²UMIC Research Centre, RWTH Aachen University, Aachen, Germany

Abstract. We propose a powerful pipeline for determining the pose of a query image relative to a point cloud reconstruction of a large scene consisting of more than one million 3D points. The key component of our approach is a new active search method that efficiently finds the correspondences between image features and scene points needed for pose estimation.

Our main contribution is a framework for actively searching for additional matches, based on both 2D-to-3D and 3D-to-2D search. A unified formulation of search in both directions allows us to exploit the distinct advantages of both strategies, while avoiding their weaknesses. Due to active search, the resulting pipeline is able to close the gap in registration performance observed between efficient search methods and approaches that are allowed to run for multiple seconds, without sacrificing run-time efficiency. Our method achieves the best registration performance published so far on three standard benchmark datasets, with run-times comparable or superior to the fastest state-of-the-art methods.

1

From Meaningful Contours to Discriminative

Object Shape

Pradeep Yarlagadda and Björn Ommer

University of Heidelberg,

Speyererstr. 6, 69115 Heidelberg, Germany

{pradeep.yarlagadda,bjoern.ommer }@iwr.uni-heidelberg.de

Abstract. Shape is a natural, highly prominent characteristic of objects that human vision utilizes everyday. But despite its expressiveness, shape poses significant challenges for category-level object detection in cluttered scenes: Object form is an emergent property that cannot be perceived locally but becomes only available once the whole object has been detected and segregated from the background. Thus we address the detection of objects and the assembling of their shape simultaneously. A dictionary of meaningful contours is obtained by clustering based on contour co-activation in all training images. We seek a joint, consistent placement of all contours in an image, since placing them independently from another is not reliable due to the emergence of shape. Therefore, the characteristic object shape is learned by discovering spatially consistent configurations of all dictionary contours using maximum margin multiple instance learning. During recognition, objects are detected and their shape is explained simultaneously by optimizing a single cost function. We demonstrate the benefit of our approach on standard shape benchmarks.

1

From Meaningful Contours to Discriminative Object Shape 779

11. Berg, A.C., Berg, T.L., Malik, J.: Shape matching and object recognition using

low distortion correspondence. In: CVPR, pp. 26–33 (2005)

12. Julesz, B.: Textons, the elements of texture perception and their interactions. Nature 29(290), 91–97 (1981)

13. Csurka, G., Dance, C.R., Fan, L., Willamowski, J., Bray, C.: Visual categorization

with bags of keypoints. In: ECCV, Workshop Stat. Learn. in Comp. Vis. (2004)

14. Gall, J., Lempitsky, V.: Class-specific hough forests for object detection. In: CVPR (2009)

15. Maji, S., Malik, J.: Object detection using a max-margin hough transform. In: CVPR (2009)

16. Felzenszwalb, P., Girshick, R., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. PAMI 32, 1627–1645 (2010)

17. Yarlagadda, P., Monroy, A., Ommer, B.: Voting by Grouping Dependent Parts. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part V. LNCS, vol. 6315, pp. 197–210. Springer, Heidelberg (2010)

18. Gavrila, D.: A bayesian, exemplar-based approach to hierarchical shape matching. PAMI 29 (2007)

19. Felzenszwalb, P., McAllester, D., Ramanan, D.: A discriminatively trained, multi-scale, deformable part model. In: CVPR (2008)

20. Toshev, A., Taskar, B., Daniilidis, K.: Object detection via boundary structure segmentation. In: CVPR, pp. 950–957 (2010)

21. Zhu, L., Chen, Y., Lin, C., Yuille, A.: Max-margin learning of hierarchical configurable deformable templates (hcdt) for efficient object parsing and pose estimation. IJCV 93, 1–21 (2011)

22. Fidler, S., Leonardis, A.: Towards scalable representations of object categories: Learning a hierarchy of parts. In: CVPR (2007)

23. Ahuja, N., Todorovic, S.: Connected segmentation tree: A joint representation of

- region layout and hierarchy. In: CVPR (2008)
24. Kokkinos, I., Yuille, A.L.: Hop: Hierarchical object parsing. In: CVPR (2009)
25. Tu, Z., Chen, X., Yuille, A., Zhu, S.: Image parsing: Unifying segmentation, detection, and recognition, vol. 2 (2005)
26. Sala, P., Dickinson, S.: Contour Grouping and Abstraction Using Simple Part Models. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part V. LNCS, vol. 6315, pp. 603–616. Springer, Heidelberg (2010)
27. Ma, T., Latecki, L.: From partial shape matching through local deformation to robust global shape similarity for object detection. In: CVPR (2011)
28. Srinivasan, P., Zhu, Q., Shi, J.: Many-to-one contour matching for describing and discriminating object shape. In: CVPR (2010)
29. Liu, M., Tuzel, O.: A. Veeraraghavan, Chellappa, R.: Fast directional chamfer matching. In: CVPR (2010)
30. Andrews, S., Tsochantaridis, I., Hofmann, T.: Support vector machines for multiple-instance learning. In: NIPS (2003)
31. Narasimhan, M., Bilmes, J.: A submodular-supermodular procedure with applications to discriminative structure learning. In: UAI, pp. 401–412 (2005)
32. Riemenschneider, H., Donoser, M., Bischof, H.: Using Partial Edge Contour Matches for Efficient Object Category Localization. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part V. LNCS, vol. 6315, pp. 29–42. Springer, Heidelberg (2010)
33. Ferrari, V., Jurie, F., Schmid, C.: From image to shape models for object detection. IJCV 87, 284–303 (2010)
34. Ommer, B., Malik, J.: Multi-scale object detection by clustering lines. In: ICCV (2009)

A Particle Filter Framework for Contour Detection

Nicolas Widynski and Max Mignotte

Department of Computer Science and Operations Research (DIRO), University of Montreal,

C.P. 6128, succ. Centre-Ville, , Montreal (Quebec), H3C 3J7, Canada

{widynski,mignotte}@iro.umontreal.ca

Abstract. We investigate the contour detection task in complex natural images. We propose a novel contour detection algorithm which locally tracks small pieces of edges called edgelets. The combination of the Bayesian modeling and the edgelets enables the use of semi-local prior information and image-dependent likelihoods. We use a mixed of offline and online learning strategy to detect the most

relevant edgelets. The detection problem is then modeled as a sequential Bayesian tracking task, estimated using a particle filtering technique. Experiments on the

Berkeley Segmentation Datasets show that the proposed Particle Filter Contour Detector method performs well compared to competing state-of-the-art methods.

1

TriCoS: A Tri-level Class-Discriminative

Co-segmentation Method for Image Classification

Yuning Chai, Esat Sabat, Victor Lempitsky,

Luc Van Gool, and Andrew Zisserman

¹Computer Vision Group, ETH Zurich, Switzerland

²Machine Vision Group, University of Oulu, Finland

³Yandex, Russia

⁴Visual Geometry Group, University of Oxford, United Kingdom

Abstract. The aim of this paper is to leverage foreground segmentation to improve

classification performance on weakly annotated datasets - those with no additional annotation other than class labels. We introduce TriCoS, a new co-segmentational algorithm that looks at all training images jointly and automatically segments out the most class-discriminative foregrounds for each image. Ultimately, those foreground

segmentations are used to train a classification system. TriCoS solves the co-segmentation problem by minimizing losses at three different levels: the category level for foreground/background consistency across images belonging to the same category, the image level for spatial continuity within each image, and the dataset level for discrimination between classes. In an extensive set of experiments, we evaluate the algorithm on three benchmark datasets: the UCSD-Caltech Birds-200-2010, the Stanford Dogs, and the Oxford Flowers 102. With the help of a modern image classifier, we show superior performance compared to previously published classification methods and other co-segmentation methods.

1

Multi-view Discriminant Analysis

Meina Kan¹, Shiguang Shan¹, Haihong Zhang²,
Shihong Lao², and Xilin Chen¹

¹Key Lab of Intelligent Information Processing of Chinese Academy of Sciences (CAS), Institute of Computing Technology, CAS, Beijing, 100190, China

²Omron Social Solutions Co., LTD., Kyoto, Japan

{meina.kan, shiguang.shan, xilin.chen}@vipl.ict.ac.cn,

angelazhang@ssb.kusatsu.omron.co.jp, lao@ari.ncl.omron.co.jp

Abstract. The same object can be observed at different viewpoints or even by different sensors, thus generating multiple distinct even heterogeneous samples.

Nowadays, more and more applications need to recognize object from distinct views. Some seminal works have been proposed for object recognition across two views and applied to multiple views in some inefficient pairwise manner. In this paper, we propose a Multi-view Discriminant Analysis (MvDA) method, which seeks for a discriminant common space by jointly learning multiple view-specific linear transforms for robust object recognition from multiple views, in an non-pairwise manner. Specifically, our MvDA is formulated to jointly solve the multiple linear transforms by optimizing a generalized Rayleigh quotient, i.e., maximizing the between-class variations and minimizing the within-class variations of the low-dimensional embeddings from both intra-view and inter-view in the common space. By reformulating this problem as a ratio trace problem, an analytical solution can be achieved by using the generalized eigenvalue decomposition. The proposed method is applied to three multi-view face recognition problems: face recognition across poses, photo-sketch face recognition, and Visual (VIS) image vs. Near Infrared (NIR) image face recognition. Evaluations are conducted respectively on Multi-PIE, CUFSF and HFB databases. Intensive experiments show that MvDA can achieve a more discriminant common space, with up to 13% improvement compared with the best known results.

Keywords: Multi-view Discriminant Analysis, Multi-view Face Recognition, Common space for Multi-view.

1

Multi-scale Patch Based Collaborative

Representation for Face Recognition

with Margin Distribution Optimization

Pengfei Zhu¹, Lei Zhang¹, Qinghua Hu², and Simon C.K. Shiu¹

¹Biometric Research Center, Dept. of Computing,

The Hong Kong Polytechnic University

²School of Computer Science and Technology, Tianjin University

{cspzhu, cslzhang}@comp.polyu.edu.hk

Abstract. Small sample size is one of the most challenging problems in face recognition due to the difficulty of sample collection in many real-world applications. By representing the query sample as a linear combination of training samples from all classes, the so-called collaborative representation based classification (CRC) shows very effective face recognition performance with low computational cost. However, the recognition rate of CRC will drop dramatically when the available training samples per subject are very limited. One intuitive solution to this problem is operating CRC on patches and combining the recognition outputs of all patches. Nonetheless, the setting of patch size is a non-trivial task. Considering the fact that patches on different scales can have complementary information for classification, we propose a multi-scale patch based CRC method, while the ensemble of multi-scale outputs is achieved by regularized margin distribution optimization. Our extensive experiments validated that the proposed method outperforms many state-of-the-art patch based face recognition algorithms.

1

Object Detection Using

Strongly-Supervised Deformable Part Models

Hossein Azizpour and Ivan Laptev

1Computer Vision and Active Perception Laboratory (CVAP), KTH, Sweden

2INRIA, WILLOW, Laboratoire d'Informatique de l'Ecole Normale Supérieure

azizpour@kth.se, ivan.laptev@inria.fr

Abstract. Deformable part-based models [1, 2] achieve state-of-the-art performance for object detection, but rely on heuristic initialization during training due to the optimization of non-convex cost function. This paper investigates limitations of such an initialization and extends earlier methods using additional supervision. We explore strong supervision in terms of annotated object parts and use it to (i) improve model initialization, (ii) optimize model structure, and (iii) handle partial occlusions. Our method is able to deal with sub-optimal and incomplete annotations of object parts and is shown to benefit from semi-supervised learning setups where part-level annotation is provided for a fraction of positive examples only. Experimental results are reported for the detection of six animal classes in PASCAL VOC 2007 and 2010 datasets. We demonstrate significant improvements in detection performance compared to the LSVM [1] and the Poselet [3] object detectors.

1

Efficient Misalignment-Robust Representation

for Real-Time Face Recognition

Meng Yang, Lei Zhang, and David Zhang

Dept. of Computing, The Hong Kong Polytechnic University, Hong Kong

{csmyang, cslzhang}@comp.polyu.edu.hk

Abstract. Sparse representation techniques for robust face recognition have been widely studied in the past several years. Recently face recognition with simultaneous misalignment, occlusion and other variations has achieved interesting results via robust alignment by sparse representation (RASR). In RASR, the best alignment of a testing sample is sought subject by subject in the database. However, such an exhaustive search strategy can make the time complexity of RASR prohibitive in large-scale face databases. In this paper, we propose a novel scheme, namely misalignment robust representation (MRR), by representing the misaligned testing sample in the transformed face space spanned by all subjects. The MRR seeks the best alignment via a two-step optimization with a coarse-to-fine search strategy, which needs only two deformation-recovery operations. Extensive experiments on representative face databases show that MRR has almost the same accuracy as RASR in various face recognition and verification tasks but it runs tens to hundreds of times faster

than RASR. The running time of MRR is less than 1 second in the large-scale Multi-PIE face database, demonstrating its great potential for real-time face recognition.

1

Monocular Object Detection Using 3D

Geometric Primitives

Peter Carr¹, Yaser Sheikh², and Iain Matthews^{1,2}

¹Disney Research, Pittsburgh

²Carnegie Mellon University

Abstract. Multiview object detection methods achieve robustness in adverse imaging conditions by exploiting projective consistency across views. In this paper, we present an algorithm that achieves performance comparable to multiview methods from a single camera by employing geometric primitives as proxies for the true 3D shape of objects, such as pedestrians or vehicles. Our key insight is that for a calibrated camera, geometric primitives produce predetermined location-specific patterns in occupancy maps. We use these to define spatially-varying kernel functions of projected shape. This leads to an analytical formation model of occupancy maps as the convolution of locations and projected shape kernels. We estimate object locations by deconvolving the occupancy map using an efficient template similarity scheme. The number of objects and their positions are determined using the mean shift algorithm. The approach is highly parallel because the occupancy probability of a particular geometric primitive at each ground location is an independent computation. The algorithm extends to multiple cameras without requiring significant bandwidth. We demonstrate comparable performance to multiview methods and show robust, real-time object detection on full-resolution HD video in a variety of challenging imaging conditions.

1
