

Vladislav Golyanik

Robust Methods for Dense Monocular Non-Rigid 3D Reconstruction and Alignment of Point Clouds

Robust Methods for Dense Monocular Non-Rigid 3D Reconstruction and Alignment of Point Clouds

Vladislav Golyanik

Robust Methods for Dense Monocular Non- Rigid 3D Reconstruction and Alignment of Point Clouds

 Springer Vieweg

Vladislav Golyanik
Computer Graphics D4
Max Planck Institute for Informatics
Saarbruecken, Germany

Vom Fachbereich Informatik der Technischen Universität Kaiserslautern zur Verleihung des akademischen Grades Doktor der Ingenieurwissenschaften (Dr.-Ing.) genehmigte Dissertation

Datum der wissenschaftlichen Aussprache: 20. November 2019
Dekan: Prof. Dr. Stefan Deßloch
Vorsitzender der Promotionskommission: Prof. Dr. Hans Hagen
Erster Berichterstatter: Prof. Dr. Didier Stricker
Zweiter Berichterstatter: Prof. Dr. Reinhard Koch
Dritter Berichterstatter: Prof. Dr. Antonio Agudo

Technische Universität Kaiserslautern, 2019
D386

ISBN 978-3-658-30566-6 ISBN 978-3-658-30567-3 (eBook)
<https://doi.org/10.1007/978-3-658-30567-3>

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Fachmedien Wiesbaden GmbH, part of Springer Nature 2020

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer Vieweg imprint is published by the registered company Springer Fachmedien Wiesbaden GmbH part of Springer Nature.

The registered company address is: Abraham-Lincoln-Str. 46, 65189 Wiesbaden, Germany

Acknowledgments

I am deeply and sincerely thankful to my doctoral supervisor Didier Stricker, the head of the Augmented Vision Department at the German Research Centre for Artificial Intelligence (DFKI) and a professor at the University of Kaiserslautern. He has always supported me in my research intentions, appreciated my independent work and strategically contributed to my development. He has also encouraged me to complete a research internship in a company.

The AV group has celebrated its tenth anniversary in August 2018. Its warm and collaborative environment has, without doubt, contributed to the accomplishment of this work. Therefore, I would like to acknowledge all colleagues and visiting researches of the AV group, who have been working in the AV group from 2014 until 2018, including my office mates, Adytia Tewari, Jilliam Maria Diaz Barros and Sarvenaz Salehi. A special thank goes to Leivy Michelly Kaul and Keonna Cunningham, for helping me to struggle through the office jungles and plenty of formalities.

I am deeply thankful to Bertram Taetz who was my postdoctoral colleague during my first months at DFKI. It was a unique experience to work with him which has effected my mindset in a way that I began to perceive computer vision in the continuous domain and from the perspective of applied mathematics. Since that time, I highly appreciate working with excellent mathematicians. I am deeply thankful to Gerd Reis who has often been helping me with improving paper drafts. Discussions with him have always been inspiring and rewarding.

During my sabbatical at NVIDIA Research in Santa Clara, I had the great fortune to work with Robert Maier who was an intern at NVIDIA Research at the same time as me, and Matthias Nießner who has been a visiting assistant professor at Stanford Univerity back then. I am thankful to my supervisors Kihwan Kim and Jan Kautz from NVIDIA Research for giving me the opportunity to work on a

challenging topic and sharing the ideas. Moreover, working with them has helped me to strengthen my leadership qualities.

I am deeply thankful to all my further co-authors: Gabriele Bleser, Torben Fetzner, Soshi Shimada, Mitra Narsi, Oliver Wasenmüller, Jan C. Peters, Aman Shankar Mathur, Sk Aziz Ali, Mohammad Dawud Ansari, Tomonari Yosida, Kiran Varanasi, André Jonas and Christian Theobalt. I have enjoyed working with all of you, and I have learned from all of you. Furthermore, I would like to thank Stella Graßhof, Hanno Ackermann, Bodo Rosenhahn, Jörn Ostermann, Antonio Agudo, Francesc Moreno-Noguer, Daniel Cremers, Willi Frieden, Gabriele Steidl, Yongzhi Su, Alain Pagani, Norbert Schmitz, Marco Paladini, Deqing Sun and Thomas Breuel for collaborations, fruitful discussions and advice.

Last but not least, I am genuinely and sincerely grateful to my parents and my family. Without their patience and support, endurance and apprehension this work would have been barely possible. At all times, the music of Ludwig van Beethoven was helping me to make decisions and served as an inexhaustible source of patience, persistence and inspiration.

Vladislav Golyanik

Contents

List of Figures	XV
List of Tables	XIX
Abstract	XXI
Zusammenfassung	XXIII
1 Introduction	1
1.1 Monocular Non-Rigid Dynamic 3D Reconstruction	2
1.2 Point Set Registration	3
1.3 Scope of the Thesis	4
1.4 Overview of the Contributions	6
1.5 Thesis Structure	9
1.5.1 Supporting Publications	10
2 Preliminaries	13
2.1 Computer Vision Primer	13
2.1.1 Perspective and Orthographic Projections	13
2.1.2 Problem Classification in the Sense of Hadamard	15
2.1.3 Inverse Problems in Computer Vision	15
2.1.4 Non-Linear Least Squares	17
2.1.4.1 Levenberg-Marquardt Algorithm	18
2.1.4.2 Huber Norm	19
2.2 From Sparse Rigid SfM to Sparse NRSfM	19
2.2.1 Eigenvalue Decomposition of a Matrix	19
2.2.2 Singular Value Decomposition	21
2.2.3 Rigid Structure from Motion by Factorisation	22
2.2.4 Non-Rigid Structure from Motion by Factorisation with Low-Rank Subspace Model	25

2.2.5	Parametrisation of Rotations	27
2.2.5.1	Axis-Angle Representation	28
2.2.5.2	Quaternions	29
2.2.6	Finding a Closest Rotation Matrix to a Given \mathbf{A}	31
2.2.7	Optical Flow Estimation	31
2.2.8	Multiframe Optical Flow with Subspace Constraints and Occlusion Handling	32
2.3	Local Refinement and Probabilistic Approaches for Point Set Registration	35
2.3.1	Estimation of an Optimal Transformation	35
2.3.2	Iterative Closest Point	36
2.3.3	N -Body Simulations	37
2.3.3.1	Acceleration Techniques for N -body Simulations	39
2.3.4	Gaussian Mixture Models and Expectation- Maximisation	39
2.3.4.1	Gaussian Mixture Models	40
2.3.4.2	Expectation-Maximisation Algorithm	41
2.3.5	Coherent Point Drift	41
3	Review of Previous Work	43
3.1	Non-Rigid Structure from Motion	43
3.1.1	Approaches to Monocular Non-Rigid Surface Recovery	45
3.1.2	Previous Work in Non-Rigid Structure from Motion	46
3.1.2.1	Dense Non-Rigid Structure from Motion	48
3.2	Point Set Registration	49
3.2.1	Scope of Point Set Registration	50
3.2.2	Previous Work in Rigid Point Set Registration	52
3.2.3	Previous Work in Non-Rigid Point Set Registration	54
4	Scalable Dense Non-Rigid Structure from Motion	57
4.1	Scalable NRSfM with Semidefinite Programming	58
4.1.1	Introduction	58
4.1.1.1	Contributions	60
4.1.2	Related Work	60
4.1.3	Accelerated Metric Projections (AMP) Approach	61
4.1.3.1	Coefficient Splitting	63
4.1.3.2	Constraints in the Unified Form	63
4.1.3.3	Constraints on the Traces	66
4.1.3.4	Constraints on the Combined Matrix \mathbf{Y}	66
4.1.4	Implementation	68
4.1.4.1	CSDP Solver	69

4.1.5	Experiments	70
4.1.5.1	Quantitative Evaluation	70
4.1.5.2	Qualitative Results on Real and Rendered Image Sequences	72
4.1.5.3	Discussion	74
4.1.6	Conclusion	74
4.2	Scalable NRSfM with Few Prior Assumptions	74
4.2.1	Introduction and an Overview of Contributions	74
4.2.2	Related Work	75
4.2.3	Scalable Monocular Surface Reconstruction Approach	75
4.2.3.1	Problem Formulation	76
4.2.3.2	Smooth Shape Trajectory	77
4.2.3.3	Non-Rigid Shape Recovery	78
4.2.4	Experimental Results	79
4.2.4.1	Dense Datasets with Ground Truth	80
4.2.4.2	Evaluation on Sparse Datasets	82
4.2.4.3	Evaluation on Dense Datasets	83
4.2.4.4	NRSfM Challenge 2017	86
4.3	Conclusion	88
5	Shape Priors in Dense Non-Rigid Structure from Motion	89
5.1	Static Shape Prior for Explicit Occlusion Handling	90
5.1.1	Motivation and Contributions	90
5.1.2	Related Work	90
5.1.3	Variational Approach with a Shape Prior (SPVA)	92
5.1.3.1	Per Sequence Shape Prior	93
5.1.3.2	Per Frame Shape Prior	97
5.1.3.3	Per Pixel Per Frame Shape Prior	97
5.1.4	Obtaining Shape Prior	98
5.1.4.1	Occlusion Tensor Estimation	98
5.1.4.2	Total Intensity Criterion	100
5.1.5	Experiments	100
5.1.5.1	Evaluation Methodology	102
5.1.5.2	Experiments on Synthetic Data	105
5.1.5.3	Experiments on Real Data	108
5.2	Intrinsic Dynamic Shape Prior for Dense NRSfM	113
5.2.1	Related Work	114
5.2.2	The Proposed Approach with Dynamic Shape Prior	115
5.2.2.1	Obtaining Dynamic Shape Prior (DSP)	117
5.2.3	Experimental Evaluation	119

5.2.3.1	Evaluation Methodology	119
5.2.3.2	Evaluation of CMDR Disjointly from DSPR	121
5.2.3.3	Self- and Cross-Convergence Tests	122
5.2.3.4	Influence of the MSGD Parameters	125
5.2.3.5	Joint Evaluation of Flow and DSPR	126
5.2.3.6	Experiments with Real Data and Applications	127
5.3	Conclusion	132
6	Coherent Depth Fields with High Dimensional Space Model	135
6.1	Depth Fields in NRSfM	136
6.1.1	Motivation and Significance of Depth Fields	136
6.1.2	Contributions	136
6.1.3	Related Work	137
6.1.4	Coherency Term	138
6.1.5	Coherent Depth Fields (CDF) Approach	139
6.1.6	Experiments	142
6.1.6.1	Synthetic Sequences and Joint Evaluation with MFOF	144
6.1.6.2	Real Sequences	146
6.2	High Dimensional Space Model for NRSfM	148
6.2.1	Motivation	148
6.2.2	Contributions	149
6.2.3	Related Work	150
6.2.3.1	Localised Modelling	150
6.2.3.2	Compressed and Compact Representations	150
6.2.3.3	Coarse-to-Fine Recovery	151
6.2.3.4	Geometry Lifting	151
6.2.4	High Dimensional Space Model (HDSM)	151
6.2.4.1	Considerations in the Rigid Case	152
6.2.4.2	Considerations in the Non-Rigid Case	153
6.2.4.3	HDSM and Other Deformation Models	154
6.2.5	Lifted Coherent Depth Fields with HDSM	155
6.2.5.1	Lifting-Compression of \mathbf{S}	156
6.2.5.2	Decompression-Expansion of the Lifted Geometry	158
6.2.5.3	Solution Initialisation	159
6.2.6	Experimental Results	159
6.2.6.1	Synthetic Face Sequences	160
6.2.6.2	Real and Naturalistic Image Sequences	161
6.3	Conclusion	162

7	Monocular Surface Regression with Learned Deformation Model	165
7.1	Architecture of the Hybrid Deformation Model Network (HDM-Net)	167
7.1.1	Loss Functions	168
7.2	Dataset and Training	173
7.3	Geometry Regression and Comparisons	173
7.4	Concluding Remarks	179
8	Probabilistic Point Set Registration with Prior Correspondences	181
8.1	Rigid Point Set Registration with Prior Correspondences	183
8.1.1	Extended Coherent Point Drift (ECPD) in General Form	185
8.1.2	Rigid Extended Coherent Point Drift (R-ECPD)	186
8.1.3	Evaluation	188
8.1.3.1	Experiments with Synthetic Data	188
8.1.3.2	Experiments with Real Data	189
8.1.3.3	The Stadium Dataset	190
8.2	Non-Rigid Point Set Registration with Prior Correspondences	192
8.2.1	Related Work	193
8.2.2	Non-Rigid Extended Coherent Point Drift (ECPD)	195
8.2.3	Implementation	197
8.2.3.1	Coarse-To-Fine Strategy with Correspondence Preserving Subsampling	198
8.2.4	Evaluation	201
8.2.4.1	Experiments with Synthetic Data	201
8.2.4.2	Experiments with Real Data	204
8.2.5	Proof of the Proposition	208
8.3	An Application in a Pipeline for Human Appearance Transfer	211
8.3.1	Related work	212
8.3.2	The Proposed Framework	213
8.3.2.1	A 3D Human Body Template	213
8.3.2.2	Overview of the Framework	214
8.3.2.3	Landmark Extraction	214
8.3.2.4	Post-Processing	216
8.3.2.5	Handling Variety in Hand Poses	216
8.3.3	Experimental Results	217
8.3.3.1	Experiments with Real Data	218
8.3.3.2	A System for Treatment of Social Pathologies	219
8.4	Summary and Conclusion	221

9	Point Set Registration Relying on Principles of Particle Dynamics	225
9.1	Rigid Gravitational Approach (GA) with Second-Order ODEs . . .	225
9.1.1	Gravitational Approach	228
9.1.1.1	Gravitational Potential Energy	230
9.1.1.2	Rigidity Constraints	231
9.1.1.3	Acceleration Techniques	233
9.1.2	Evaluation	235
9.1.2.1	Experiments on Synthetic Data	235
9.1.2.2	Experiments on Real Data	238
9.1.2.3	Experiments on SLAM Datasets	239
9.1.2.4	Discussion	245
9.2	Accelerated Gravitational Approach with Altered Laws of Physics	245
9.2.1	The Enhanced Particle Dynamics Based Gravitational Approach	247
9.2.1.1	Acceleration with a Barnes-Hut Tree	247
9.2.1.2	Local Enhancement with Spherical Coordinates	250
9.2.1.3	Handling Varying Point Densities	251
9.2.1.4	Energy Minimisation	251
9.2.2	Experimental Evaluation	252
9.2.2.1	Quantitative Evaluation	252
9.2.2.2	Evaluation with Real-World Data	259
9.3	Gravitational Approach for Non-Rigid Point Set Registration . . .	260
9.3.1	Non-Rigid Gravitational Approach (NRGA)	261
9.3.1.1	Modified N -Body Problem	263
9.3.1.2	Distributed Locally Multiply-Linked Policy . . .	263
9.3.1.3	Coherent Collective Motion Regulariser	265
9.3.1.4	Algorithm and Complexity Analysis	266
9.3.2	Experimental Evaluation	266
9.3.2.1	Evaluation Methodology and Datasets	267
9.3.2.2	Experimental Results on Synthetic Data	269
9.3.2.3	Experimental Results with Qualitative Interpretation	270
9.4	Conclusion	272
10	Applications to Scene Flow Estimation	275
10.1	Scene Flow from Monocular Image Sequences	275
10.1.1	Monocular Scene Flow as an Emerging Field	276
10.1.2	MSF and NRSfM in the Continuous Domain	279
10.1.3	The NRSfM-Flow Framework	283

10.1.4	Evaluation	285
10.1.5	Conclusion	289
10.2	RGB-D Multiframe Scene Flow with Piecewise Rigid Motion . . .	290
10.2.1	Motivation, Preliminaries and Contributions	290
10.2.2	Previous and Related Works in the Area of RGB-D Scene Flow Estimation	292
10.2.3	Multiframe Scene Flow (MSF) with Piecewise Rigid Moti- on	294
10.2.3.1	Multiframe Formulation	299
10.2.4	Energy Optimisation	301
10.2.4.1	Energy Initialisation and Settings	303
10.2.5	Experimental Evaluation	304
10.2.5.1	Experiments on Synthetic Data	305
10.2.5.2	Experiments on Real Data	308
10.2.6	Discussion	308
10.2.7	Conclusion	312
11	Summary, Conclusions and Outlook	313
11.1	Future Directions	315
	Bibliography	317
	Publication List	349

List of Figures

2.1	Energy function of the multiframe optical flow approach	33
2.2	NRSfM reconstructions of a human heart	34
3.1	Depth ambiguity in NRSfM	44
3.2	Motion and deformation cues in NRSfM	48
4.1	Runtimes of AMP on the synthetic flag sequence	64
4.2	Qualitative evaluation of AMP	69
4.3	Results of AMP on the face sequence	72
4.4	Results of AMP on the heart sequence	73
4.5	Comparison of the normalised mean 3D error (log scale)	80
4.6	Qualitative evaluation of SMSR	81
4.7	Qualitative comparison of MP, CSF1, CSF2, PTA, VA and our SMSR (<i>Actor1 Sparse</i>)	83
4.8	Qualitative evaluation of MP, PTA, CSF1, VA and our SMSR (<i>synthetic faces</i>)	84
4.9	Visualisation of the 3D motion fields recovered by SMSR	85
4.10	Qualitative evaluation of SMSR	87
5.1	An overview of the SPVA pipeline	92
5.2	Exemplary frames from the modified flag sequences	101
5.3	Plots of the total intensity function	102
5.4	Quantitative evaluation of MFSF + SPVA in different modes	103
5.5	Quantitative evaluation on the flag sequence	103
5.6	Rigid initialisation and the shape prior (overlay)	104
5.7	Exemplary frames of the <i>hashtag</i> sequence	104
5.8	Qualitative results of SPVA in comparison to other pipeline combinations	107
5.9	Experimental results on the <i>heart</i> sequence	109
5.10	Experimental results on the <i>face</i> sequence	110

5.11	Experimental results on the <i>ASL</i> sequence	111
5.12	Results on the <i>ASL</i> sequence with correspondence correction . .	112
5.13	Results of the <i>self-convergence</i> and <i>cross-convergence</i> tests . . .	120
5.14	Reconstruction of SMSR on the perturbed point tracks	122
5.15	Results of the experiments with MSGD parameters	124
5.16	Convergence patterns observed in the self- and cross-convergence tests	125
5.17	The new <i>actor mocap</i> sequence	127
5.18	A real image sequence and non-rigid 3D reconstructions thereof .	128
5.19	Exemplary reconstructions of CDF and DSPR on noisy point tracks	129
5.20	Application of DSPR in heart bypass surgery with reoccurring deformations	130
6.1	Explanation of the coherency term	138
6.2	Exemplary reconstructions by VA, AMP and our CDF on the new sequence	143
6.3	Evolution of reconstructed occluded regions for different σ . . .	144
6.4	CDF reconstruction of the laparoscopic sequence	146
6.5	Examples of shaded surfaces reconstructed by CDF	147
6.6	An overview of the main idea	149
6.7	e_{3D} and ϵ as functions of ϵ	158
6.8	Visualisations of final permutation matrices Π and series of Φ . .	159
6.9	Reconstructions of several frames of the <i>synthetic face</i> by VA, MP, TB and L-CDF	160
6.10	Exemplary reconstructions of real and synthetic image sequences	163
7.1	Reconstruction of an endoscopically textured surface	166
7.2	An overview of the HDM-Net architecture	169
7.3	Encoder and decoder of HDM-Net	170
7.4	Contour loss	171
7.5	Camera poses used for the dataset generation	171
7.6	The pattern of the training and test datasets	173
7.7	Selected reconstruction results on endoscopically textured surfaces	174
7.8	The effect of the isometry prior	175
7.9	Qualitative comparisons of the results of HDM-Net, AMP, VA and Yu <i>et al.</i>	177
7.10	Exemplary reconstructions of HDM-Net on real images	178
7.11	Graphs of e_{3D} for HDM-Net	178

8.1	Embedding of prior correspondences into probabilistic point set registration	182
8.2	Rigid point cloud registration with prior matches	184
8.3	Rigid registration of <i>Lion</i> dataset	190
8.4	Registration of partially overlapping shapes	191
8.5	Non-rigid registration of point sets representing arms in different poses	193
8.6	Non-rigid registration of a 2D “Fish” dataset	199
8.7	Acceleration scheme of ECPD	200
8.8	Results of the ECPD experiment with SINTEL dataset	203
8.9	Comparison of the registration results of CPD and ECPD	205
8.10	Registration of the “woman with a scarf” dataset	205
8.11	Non-rigid registration of the “man with a hood” dataset	206
8.12	An overview of the human appearance transfer framework	211
8.13	A full-body 3D human template	213
8.14	Extraction of the body landmarks	215
8.15	Avoiding hand flattening	216
8.16	Accuracy evaluation of the proposed approach	218
8.17	Results on the FAUST dataset	220
8.18	Registration result (a template with 10^4 points)	220
8.19	The proposed framework in the treatment of social pathologies	222
9.1	Point set registration with the gravitational approach	226
9.2	Registration results of ICP, CPD and our GA on data with clustered outliers	227
9.3	Registration results on <i>Stanford bunny</i>	236
9.4	Registration results on data with uniformly distributed and Gaussian noise	237
9.5	Results on data with structured outliers and missing parts	238
9.6	Experiment with prior correspondences as applied to image registration	240
9.7	Selected point clouds converted from the depth maps	241
9.8	Results of the experiment on the Stanford 3D datasets	242
9.9	Depth maps involved in the experiment	243
9.10	Results of the experiment on the CoRBS dataset	244
9.11	An overview of the proposed BH-RGA approach	246
9.12	Clusters fetched during alignment (<i>clean-500</i> experiment)	249
9.13	RMSE as a function of the point perturbation magnitude index (<i>U256</i> and <i>G256</i>)	255
9.14	Fragment comparison after the alignment (<i>sleeping2</i>)	256

9.15	Runtime evaluation metrics as the functions of the threshold γ . . .	257
9.16	Examples of reprojected 3D flows obtained by BH-RGA	257
9.17	Alignment of partially overlapping real-world data	259
9.18	Usage of the law of universal gravitation for non-rigid point set alignment	261
9.19	Two main steps of NRGGA	262
9.20	Visualisation of point trajectories during alignment, with and without CCM regulariser	265
9.21	Quantitative results of NRGGA on several datasets	268
9.22	Handling of missing data	269
9.23	Qualitative results of NRGGA, CPD, GMMReg and NR-ICP on human face scans	271
9.24	The experiment with real data with per-point distance error . . .	272
10.1	Results of NRSfM-Flow on the <i>human face</i> sequence	276
10.2	Overview of NRSfM-Flow	280
10.3	Geometric interpretations	282
10.4	Experimental results on the <i>barn owl</i> sequence	286
10.5	Examples of Poisson reconstructions	286
10.6	Experimental results on the SINTEL dataset	287
10.7	Results on the <i>heart</i> and <i>music notes</i> sequences	288
10.8	A high-level overview of the proposed MSF approach	291
10.9	An overview of the main related works	293
10.10	An overview of the main components of the proposed energy . . .	296
10.11	Projective ICP term	297
10.12	Segmentation transfer from the reference frame to three other frames (<i>alley1</i>)	300
10.13	Experimental results on the SINTEL <i>alley1</i> and <i>bandage1</i>	302
10.14	Experimental results on a static scene observed by a moving camera (SINTEL <i>sleeping2</i>)	303
10.15	Results on the Bonn multibody dataset	304
10.16	Results on several real datasets (<i>Chairs</i> , <i>Pile of Boxes</i> and <i>Board</i>)	306
10.17	Segmentation transfer on the Bonn watering can sequence	307
10.18	The order of frames used in the remaining figures	308
10.19	Additional visualisations of the <i>alley1</i> sequence	309
10.20	Additional visualisations of the <i>bandage1</i> , <i>sleeping1</i> and <i>shaman2</i> sequences	310
10.21	Additional visualisations of the <i>mountain1</i> , <i>sleeping2</i> and <i>shaman3</i> sequences	311

List of Tables

1.1	List of supporting publications	10
1.2	List of conference proceedings with abbreviations	11
4.1	e_{3D} for benchmark datasets (sparse)	82
4.2	e_{3D} for benchmark datasets (dense)	84
4.3	e_{3D} for the modified benchmark dataset	85
5.1	RMSE of different algorithmic combinations	108
5.2	Runtimes of different algorithm combinations	112
5.3	Parameters of the proposed approach	112
5.4	RMSE of several methods (<i>synthetic faces</i>)	120
5.5	Quantitative comparison of CMDR to several other methods . . .	121
5.6	Compression ratios	131
6.1	RMSE of VA, AMP and the proposed CDF on the <i>actor</i> dataset . .	142
6.2	Average RMSE on the occluded flag sequences	145
6.3	Joint average RMSE and s on the synthetic faces	146
6.4	A non-exhaustive list of symbols used in the section	152
6.5	Joint average e_{3D} and σ_e for the <i>synthetic faces</i>	160
7.1	Comparisons of per-frame runtime t , e_{3D} and σ	176
7.2	Comparison of 3D error for different textures	176
7.3	Comparison of 3D error for different illuminations	176
7.4	Comparison of effects of loss functions	177
8.1	Speedup of ECPD (“woman with a scarf”)	207
8.2	Speedup of ECPD (“man with a hood”)	208
8.3	The parameters for the core steps of the proposed pipeline	217
9.1	Summary of the qualitative evaluation of the compared methods . .	253
9.2	RMSE and σ in $U256$ and $G256$ experiments	254

9.3 Comparison of E-CPD [123] and BH-RGA with prior matches. . . 254

10.1 Core equations of the framework relating NRSfM and MSF 282

10.2 Comparison between scene flow projections and the ground truth
optical flow 306

10.3 Comparison of average runtimes of MSF and SRSF 306

Abstract

An accurate acquisition and processing of 3D point cloud data is an active research area in computer vision encompassing various unsolved problems. The thesis at hand addresses the jointly studied domains of dense non-rigid 3D reconstruction from monocular image sequences and point set registration under rigid as well as non-rigid transformations. Monocular non-rigid 3D reconstruction, which is in the focus of this work — known as non-rigid structure from motion (NRSfM) — relies on weak assumptions about the feasible deformation modes imposed on top of the motion and deformation cues. NRSfM and non-rigid point set registration are highly ill-posed problems in the sense of Hadamard.

The proposed dense NRSfM methods address the broad range of research questions including occlusion handling, scalability, interactive yet accurate processing as well as dense structure compression. For the occlusion handling and dealing with inaccurate point tracks, we propose a shape prior obtained on-the-fly and a new spatial regulariser — the coherency term. We also introduce a new model for NRSfM, which allows representing the recovered structure compactly.

The proposed point set registration methods aim at the enhanced registration accuracy for noisy data and samples with clustered outliers. For that reason, we embed prior correspondences into probabilistic point set registration and introduce a previously unexplored class of methods relying on principles of particle dynamics with simulated gravitational forces.

The thorough experimental evaluation confirms the efficiency and high accuracy of the proposed methods as well as the validity of the new ideas. By using the new principles, we advance the state of the art in dense monocular non-rigid 3D reconstruction and alignment of noisy point sets. Applications of the proposed NRSfM methods include (but are not limited to) 3D recovery and analysis of human and animal faces, endoscopic scenes and various other deformable surfaces. The proposed point set registration methods can be applied in robotics, automotive

driving, face and shape recognition, and other areas. Apart from the abovementioned applications, we show how both method classes can be used for human appearance transfer, multiframe scene flow estimation from RGB-D as well as monocular image sequences. The developed methods offer numerous avenues for further investigation.

Zusammenfassung

Die genaue Eingabe und Verarbeitung von 3D Punktwolken ist ein aktives Forschungsfeld im maschinellen Sehen, das viele ungelöste Probleme umfasst. Die vorliegende Doktorarbeit befasst sich mit den in Zusammenarbeit erforschten Bereichen der dichten nicht-starren 3D Rekonstruktion aus monokularen Bildsequenzen, mit sowohl starren als auch nicht-starren Punktwolkenregistrierung. Die monokulare 3D Rekonstruktion, die im Fokus dieser Arbeit steht und die als nicht-rigide Struktur aus Bewegung (NRSaB) bekannt ist, wertet, einerseits, Bewegungen und Deformationen aus, und, andererseits, verknüpft diese mit den zusätzlichen Annahmen und Vorwissen über die Szene und die Art der zulässigen Zustände.

Die eingeführten Verfahren zur dichten NRSaB gehen auf mehrere offene Fragen ein, und zwar auf die Behandlung von Verdeckungen, die Skalierbarkeit und die Anpassungsfähigkeit auf unterschiedliche Szenarien und Größenordnungen der Szenen, interaktive und präzise Verarbeitung, sowie die Komprimierung dichter 3D Geometrie. Zwecks der Behandlung von Verdeckungen und fehlerhafter Punktkorrespondenzen werden Verfahren mit dem am Anfang einer Bildsequenz gewonnenen Formvorwissen sowie einem neuen räumlichen Kohärenz Regularisierer vorgestellt. Darüber hinaus, leiten wir ein neues NRSaB Verfahren her, das die gewonnene Geometrie in eine kompakte Repräsentation überführt.

Die entwickelten Verfahren zur Punktwolkenregistrierung verfolgen das Ziel, verrauschte und partielle Eingabedaten mit höherer Präzision zu verarbeiten als die Vorgängermethoden. Dementsprechend schlagen wir vor, die im Vorfeld hergestellten Korrespondenzen ins probabilistische Framework für die Punktwolkenregistrierung zu integrieren, und, zweitens, präsentieren wir eine neue und bisher unerforschte Verfahrensklasse, welche die Teilchenbewegungen unter virtuellen Schwerkraften simuliert.

Durch gründliche und zahlreiche Experimente ist es uns gelungen, die Geltung der neuen Ideen sowie die Präzision und Robustheit der entwickelten Verfahren

zu bestätigen. Dank der neuen Prinzipien und Verfahren waren wir imstande, den Stand der Technik in beiden Bereichen der monokularen nicht-rigiden 3D Rekonstruktion sowie Punktwolkenregistrierung zu verbessern. Zu den Anwendungen neuer NRSaB Verfahren zählen die 3D Rekonstruktion von Menschen, Tieren und endoskopischer Aufnahmen sowie die Erfassung dünner Strukturen unterschiedlicher Herkunft. Die entwickelten Verfahren zur Punktwolkenregistrierung können unter anderem in Robotik, selbstfahrender Fahrzeugtechnik sowie Gesichts- und Formerkennung angewendet werden. Neben der erwähnten Gebieten wird in dieser Arbeit gezeigt, wie die beiden Verfahrensklassen zwecks der Übertragung des äußeren Erscheinungsbildes von Menschen sowie der Berechnung vom Szenenfluss aus Tiefenkamerabildern und monokularen Bildern angepasst werden können. Ferner, bieten die entwickelten Verfahren verschiedene Wege und Möglichkeiten zur Verbesserung und Weiterentwicklung, auf die am Schluss eingegangen wird.



1 Introduction

ONE of the objectives of computer vision is an accurate sensing of the real world and robust processing of the acquired data. Along with the material properties, knowledge of 3D geometry is a key component of complete scene description. 3D machine perception is the foundation for multiple applications which involve scene replication, scene understanding, localisation as well as a realistic superimposition of virtual contents, among others.

There are multiple sensors which come into question while designing a vision-based system including time-of-flight cameras, stereo cameras, lidars and sonars. A lightweight alternative to those is a single monocular camera. The advantages of a monocular camera are different designs and form-factors, affordability, relatively low electric energy consumption but also pervasiveness in modern electronic devices and wide acceptance in society. There are monocular cameras embedded in augmented reality glasses, helmets, mobile phones and tablets. Monocular cameras are central components in endoscopic surgery systems, surveillance systems, person identification systems, unmanned aerial and underwater vehicles, mobile robots, rovers for planetary explorations and autonomous cars. Thus, methods using monocular cameras for 3D sensing are of high relevance in a broad variety of systems and applications. Moreover, techniques for processing and analysis of the recovered raw 3D representations — often point sets and point clouds — are increasingly gaining relevance.

3D reconstruction is an extensively studied inverse problem in computer vision consisting in the recovery of the depth dimension of a scene lost during the imaging (together, the scene geometry), from single or multiple views. *Point set registration* is a computer vision problem of recovery the transformations aligning one or multiple point sets (raw 3D representations or 3D reconstructions) into a common coordinate frame or deforming the point sets so that their appearances match.

Depending on the available input and in many practical situations, 3D reconstruction can also be an ill-posed problem in the sense of Hadamard. Thus, 3D reconstruction from a single image is ill-posed, as multiple 3D scenes can result in the same 2D image. Additional prior knowledge is required to disambiguate the reconstruction such as a known object class, symmetry prior or a geometric prior. Starting from two views, additional constraints can be used ranging from epipolar geometry prior and trilinear constraints to the consistency constraints over multiple

views. Moreover, the rigidity assumption disambiguates the problem well, and impressive results were achieved in 3D reconstruction under the rigidity assumption, both from multiple views and the sequences of monocular views.

If several different views of the same scene at the same time frame are available, the subclass of the techniques is referred to as *multiple view geometry reconstruction*. In contrast, if multiple monocular views of the scene over several time frames are available, the subclass of the techniques is called *structure from motion* (SfM). In a general context, the input data corresponds to an image sequence, and some methods operate on a set of tracked points over the input views. The difference between multiple view geometry and structure from motion under rigidity is often subtle. In many cases, the techniques can be applied interchangeably, though the information about whether the cameras are static or moving can be advantageous (*e.g.*, motion blur prone to a moving camera can be accounted for). Compared to multi-view reconstruction, SfM often assumes smaller frame-to-frame displacements, as those observed in a video sequence. Video sequences also allow for stronger priors such as temporal smoothness.

Apart from predominantly static environments and sceneries, our surroundings are inhabited by living species including ourselves which move and deform. Besides, there are rigidly moving manmade instruments and products violating the assumption about the static and conserved ambience. Thus, capturing and processing of dynamic scenes is a core capability of robust vision-based systems.

1.1 Monocular Non-Rigid Dynamic 3D Reconstruction

The situation changes considerably, if the rigidity assumption does not hold anymore, *i.e.*, the scene undergoes non-rigid deformations. In the case of multiple views, the observations are captured at the same moment of time, and the geometry is still related by spatial rigidity between the views. If captured at different time frames and with different camera poses, the scene is observed in different states and the temporal rigidity does not hold anymore. The class of methods assuming non-rigid scenes over a temporal sequence of views is specified as *non-rigid structure from motion* (NRSfM). In NRSfM, the camera is moving, whereas the scene is moving and deforming. Similar to rigid SfM, the input of NRSfM is a set of tracked points over the available views.

Though NRSfM is a highly ill-posed problem which is sometimes said to be equivalent to the reconstruction from a single view, additional constraints can help to disambiguate it. *Real-world objects do not deform arbitrarily and rather follow a*

certain deformation pattern. The deformation pattern is also often associated with periodicity, which implies that the scene states are repeated in a temporally-disjoint manner. Moreover, an average or middle state can be distinguished among all observed states. Additionally, it is more probable that the states in neighbouring frames are more similar than states in frames that are temporally further apart.

There is a substantial difference between single view rigid reconstruction and NRSfM which has crystallised out. In the single view rigid reconstruction, it is valid and common to assume a specific object class, and supervised learning methods are often applied. NRSfM, in contrast, assumes that no prior shape information about the observed scene is available, and solely relies on motion and deformation cues to obtain 3D surface reconstructions from monocular image sequences. This makes NRSfM capable of handling equally well — depending on the accuracy of correspondences — thin surfaces of different kinds (flags, sails, *etc.*), human and animal faces, clothes and body tissues in medical contexts.

Several new method classes have emerged which constrain the context of NRSfM, such as those assuming an accurate reconstruction of at least one of the frames in the sequence (template-based methods) and those assuming a pre-defined deformation model but different material properties.

1.2 Point Set Registration

When 3D surface recovery is complete, there are multiple ways how the dynamic reconstruction can be processed and analysed. One of the essential pre-processing steps is changing the reference frame or pose of the reconstruction for the further comparison, deformation transfer or recognition. This operation can be performed with *rigid point set registration* if the orientation of the reference frame or object is known. The comparison and deformation transfer can be accomplished with *non-rigid point set registration*.

The objective of point set registration is to align two or several point sets, *i.e.*, to recover a transformation which registers a *template* point set to an unaltered *reference*. A point set is an unordered set of coordinates (2D or 3D), with no further information available. As a representation of a shape, it can contain noise and clustered outliers, and some parts can be missing. Point set registration should not be confused with mesh registration methods (meshes are more complete shape representations consisting of points, triangles, normals *etc.*). 3D reconstructions obtained with NRSfM often represent point sets.

In the rigid case, the transformation is parametrised by the variables of rigid body motion with six degrees of freedom (three for rotation and three for translation). During a rigid transformation, no deformation is happening, and all distances

between the points are preserved. Transformation of every point is given by the same rotation and translation. In the non-rigid case, due to deformations, distances between the points are not preserved, and the transformation is described by a general per-point displacement field. As monocular deformable reconstruction, *non-rigid point set registration relies on prior knowledge that real-world objects and scenes do not deform arbitrarily but rather follow some deformation rules and patterns*. One of the most commonly used and reasonable constraints in non-rigid point set registration is the topology preserving constraint which states that point topology must be preserved despite the distances between the points are changing. It prevents intersections between the displacements, and, as a consequence, self-intersections of the surfaces represented by the points sets (though, point sets can also represent volumetric structures). Similarly to NRSfM, non-rigid point set registration is an ill-posed problem in the sense of Hadamard.

Despite the progress in point set registration which enabled various practical applications, one of the central research questions in point set registration remains improvement of the robustness to noise and disturbing effects in the data (missing parts and clustered outliers). Moreover, processing of large point sets is an ever-relevant problem (in other words, point sets containing one-two orders of magnitude more points than what is considered as a standard nowadays; the contemporary standard in NRSfM is around $30k$ points). It is addressed with faster hardware, parallelisation as well as data structures for acceleration. In contrast to methods for processing of synthetic 3D data (computer graphics), methods for processing of raw sensor inputs have to cope and consider noise and incompleteness of the data.

3D reconstruction and point set registration exhibit similarities. Thus, common types of assumptions and constraints can be applied to disambiguate them (*e.g.*, rigidity assumption and shape priors), and for handling non-rigid deformations, regularisation of displacement fields is required. In this thesis, the study of 3D reconstruction and point set registration is conducted jointly. As will be shown throughout the thesis, both related research fields facilitate and enrich each other with ideas. Point set registration provides tools for 3D reconstruction (algorithmic and evaluation tools), 3D reconstruction provides data for optimal evaluation of point set registration, and multiple concepts can be borrowed from one field to another one (regularisation of the displacement fields).

1.3 Scope of the Thesis

This thesis focuses on robust methods for dense monocular non-rigid 3D reconstruction from uncalibrated views and alignment of point cloud data.

The methods for dense monocular non-rigid 3D reconstruction should not assume that the calibration is known, though if it is known, the algorithms could optionally use the calibration parameters. Moreover, the new approaches should reconstruct the scene per-point densely, and, optionally, allow sparse reconstruction. The requirements to the new methods include robustness to self- and external occlusions, scalability, higher accuracy and lower runtime compared to the existing methods. Some of the requirements are not necessary facilitating towards the other ones, *i.e.*, it is more challenging to develop a scalable, accurate and fast method at the same time. Efficient NRSfM methods in conjunction with robust methods for dense correspondences would enable new applications based on commodity hardware.

Thanks to the point set registration, the reconstructed scenes can be compared to some reference data or the recovered appearance can be transferred to some other representations usable in different application scenarios. Thus, both method classes can be used in a single 3D reconstruction and processing pipeline.

There is also another reason to study the fields of monocular 3D reconstruction and point cloud alignment jointly. Even though the underlying methods pursue different goals and assume different input data, both fields are still related to each other, so that cross-fertilisation and exploitation of synergies is possible. Thus, non-rigid registration can help in the joint evaluation of NRSfM and correspondence establishment approaches, as will be shown in §5. Moreover, due to the handling of deformable structures in both algorithm classes, we proposed a new spatial regulariser (coherency constraint, §6) for NRSfM which was previously used exclusively in non-rigid point set registration.

The work at hand was also inspired by the maturing research area of augmented reality. Augmented reality is an interdisciplinary research field on the intersection of computer vision, computer graphics and hardware systems (which include material science, physics, mechatronics and electronics). The goal of augmented reality is to extend and enhance the perceived reality through useful virtual contents. Virtual contents should be realistic and indistinguishable from the real ones. Along with the realistic rendering, accurate placement of virtual contents is one of the quality factors. The acquisition of geometry of deformable objects with efficient methods for processing of the reconstructions is highly relevant for augmented reality as well. Both method classes addressed in this thesis — NRSfM and point set registration — can be used in augmented reality systems in a pipeline for 3D reconstruction and data processing with a single monocular moving camera.

1.4 Overview of the Contributions

The primary subject of the dissertation is dynamic 3D reconstruction of non-rigidly deforming scenes from monocular image sequences as well as processing of point sets. The main considered algorithm classes are non-rigid structure from motion (NRSfM) and point set registration (PSR). NRSfM is a highly ill-posed inverse problem. The input of NRSfM is a set of point tracks over several unsynchronised and uncalibrated views, and the objective is the recovery of the observed non-rigid 3D geometry. Thus, NRSfM uses motion and deformation cues as well as additional weak prior assumptions about the type of valid deformations for 3D recovery. In PSR, the inputs are two point sets with a different number of points, and the objective is the alignment of those into a common reference frame (in the rigid case) or the recovery of the displacements and correspondences non-rigidly aligning the inputs (in the non-rigid case). The two fields were studied jointly and complemented each other.

NRSfM and Monocular Surface Recovery

In the field of NRSfM, the thesis features the following contributions:

- First, a new dense variational NRSfM technique for handling large occlusions and inaccuracies in the data was proposed — *Shape Prior Variational Approach* (SPVA). SPVA estimates a shape prior from several first unoccluded frames of the sequence on-the-fly and guides the reconstruction by the occlusion tensor. The occlusion tensor is computed from the initial dense flow fields and indicates occlusion probabilities for every frame. The method allows for the reconstruction of scenes where large occlusions are expected (*e.g.*, in medical contexts). The method is parallelisable and is implemented on a GPU. The experimental results show the state-of-the-art accuracy on challenging sequences for which a shape prior can be obtained.
- Second, a new method with an intrinsic dynamic shape prior for 3D reconstruction and compression of sequences with temporally-disjoint rigidity is introduced. Temporally-disjoint rigidity occurs in most real video sequences, *i.e.*, the phenomenon of state reoccurrence. The repeating states can be separated by an arbitrary number of other states and can reappear in different poses. Our *Dynamic Shape Prior Reconstruction* (DSPR) approach takes advantage of temporally-disjoint rigidity and allows for dense reconstructions with low latencies. Experiments demonstrate that DSPR can operate on inaccurate correspondences.

- Third, a new spatial regulariser — the *coherency term* — for dense NRSfM is proposed which has allowed handling of large occlusions without a shape prior. The coherency term was adopted from the motion coherence theory. Before, the coherency term was used in non-rigid point set registration. We have shown how to minimise energy with the coherency term in the context of NRSfM.
- Fourth, we have addressed the problem of structure compressibility in the sense of data compression theory in the context of NRSfM and proposed a new *High-Dimensional Space Model* (HDSM) for NRSfM. In HDSM, non-rigid geometry in 3D is encoded as multiple projections of a single high dimensional structure onto different 3D subspaces. The proposed representation in combination with the factorisation-based (decoupled) formulation for camera pose and shape recovery allows compressing the structure in the high dimensional space. The resulting method encompassing handling of inaccurate point tracks with the coherency term and structure compression is known as *Lifted Coherent Depth Fields* (L-CDF).
- Fifth, we propose a new fast technique for dense NRSfM — *Accelerated Metric Projections* (AMP) — which allows to factorise dense batches of point tracks in seconds on a CPU. At the moment of publication, AMP was the fastest dense NRSfM method delivering high reconstruction accuracy. We have shown in AMP how to minimise a quadratic function on a set of orthonormal matrices using an efficient semidefinite programming solver. The method allows an arbitrary reshuffling of the per-frame measurements which can be advantageous in the cases when temporal information cannot be maintained.
- Sixth, we have addressed the question of scalability in the context of NRSfM. The core characteristic of the resulting robust NRSfM technique — *Scalable Monocular Surface Reconstruction* (SMSR) — is the steady high accuracy across a large variety of dense and sparse datasets with reasonable runtime and linear scalability w.r.t. the number of points. In SMSR, the camera pose is updated with singular value thresholding and proximal gradient techniques, whereas the surface is estimated by alternating direction method of multipliers.
- Seventh, we found a new way to regress non-rigid geometry with a trained encoder-decoder deep neural network. In the *Hybrid Deformation Model Network* (HDM-Net), the deformation model is learned from synthetic data in a supervised manner. Among contributions of HDM-Net is a new way to perform a convolution on a point set instead of a volumetric representation, an isometric loss and a contour loss. Moreover, the inference of a surface with over 5k points takes around 5 ms. Results on real images demonstrated the potential of the proposed architecture for augmented reality applications.