

# **Efficient & Effective Image-Based Localization**

Von der Fakultät für Mathematik, Informatik und Naturwissenschaften der  
RWTH Aachen University zur Erlangung des akademischen Grades eines  
Doktors der Naturwissenschaften genehmigte Dissertation

vorgelegt von Diplom-Informatiker

**Torsten Sattler**

aus Düsseldorf, Deutschland

Berichter: Prof. Dr. Leif Kobbelt  
Prof. Dr. Bastian Leibe  
Prof. Dr. Marc Pollefeys

Tag der mündlichen Prüfung: 25.10.2013

Diese Dissertation ist auf den Internetseiten der Hochschulbibliothek online verfügbar.



# **Selected Topics in Computer Graphics**

herausgegeben von  
Prof. Dr. Leif Kobbelt  
Lehrstuhl für Informatik 8  
Computergraphik & Multimedia  
RWTH Aachen University

Band 11

**Torsten Sattler**

**Efficient & Effective Image-Based Localization**

Shaker Verlag  
Aachen 2014

**Bibliographic information published by the Deutsche Nationalbibliothek**

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <http://dnb.d-nb.de>.

Zugl.: D 82 (Diss. RWTH Aachen University, 2013)

Copyright Shaker Verlag 2014

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publishers.

Printed in Germany.

ISBN 978-3-8440-2743-3

ISSN 1861-2660

Shaker Verlag GmbH • P.O. BOX 101818 • D-52018 Aachen

Phone: 0049/2407/9596-0 • Telefax: 0049/2407/9596-9

Internet: [www.shaker.de](http://www.shaker.de) • e-mail: [info@shaker.de](mailto:info@shaker.de)

# Abstract

The problem of image-based localization is the problem of accurately determining the position and orientation from which a novel photo was taken relative to a 3D representation of the scene. It is encountered in many interesting applications such as pedestrian or robot navigation, Augmented Reality, or Structure-from-Motion, creating a strong need for algorithms solving the image-based localization problem. In this thesis, we therefore present solutions to this problem that are both effective and efficient, *i.e.*, we propose methods that can localize novel query images taken under a wide range of viewing conditions while requiring only a small amount of processing time.

We assume that the 3D scene representation is obtained by using Structure-from-Motion techniques to reconstruct the environment from a set of photos. As a result, we can associate each 3D point with multiple image descriptors modeling the local appearance of the scene around this point. We can then obtain 2D-3D correspondences between 2D feature points in the query image and 3D scene points in the model by solving a descriptor matching problem. These 2D-3D matches can in turn be used to estimate the camera position of the query image, *i.e.*, the position and orientation from which it was taken. The main difficulty of descriptor matching lies in the sheer size of the problem, since our models contain millions of 3D points while thousands of features are found in our query images. As a major contribution, we show that the resulting descriptor matching problem can still be solved very efficiently using prioritized search. We propose a prioritization scheme that is easy to implement, yet can be expected to perform close to optimal in practice. By combining our prioritization with a novel active search step that is able to discover additional matches, we are able to derive an image-based localization approach that achieves or surpasses state-of-the-art effectiveness while offering the fastest run-times published so far.

Analyzing such direct matching methods, we demonstrate that their major advantage, namely their ability to identify a set of high-quality matches, also prevents their scalability to larger datasets. Consequently, we also consider image retrieval methods for image-based localization since they are inherently more scalable. As a second major contribution, we identify the algorithmic factors preventing image retrieval methods to achieve the same effectiveness as our original system and propose a modification that is able to close the gap in effectiveness without sacrificing scalability.



# Zusammenfassung

Das Ziel von bildbasierten Lokalisierungsverfahren ist es, für ein gegebenes Fotos die Position und Ausrichtung der dazugehörigen Kamera relativ zu einem 3D Szenenmodell zu bestimmen. Das entsprechende Problem der bildbasierten Lokalisierung findet dabei viele praktische Anwendungen, wie z.B. Fußgängernavigation, Augmented Reality und Structure-from-Motion. In dieser Arbeit stellen wir effektive und effiziente Ansätze zur Lösung dieses Problems vor, d.h., wir präsentieren Verfahren welche die Position und Orientierung der Kamera für eine große Bandbreite von Blickpunkten und Beleuchtungsbedingungen in kurzer Zeit berechnen können.

Im folgenden gehen wir davon aus, dass das 3D Szenenmodell durch eine Structure-from-Motion Rekonstruktion der Umgebung aus einer Menge von Bildern erzeugt wurde. Dies erlaubt es uns jedem 3D Punkt mehrere Featuredeskriptoren zuzuweisen, welche das Aussehen der Szene um diesen Punkt herum beschreiben. Folglich können wir 2D-3D Korrespondenzen zwischen Featurepunkten im Anfragebild und 3D Punkten im Modell mit Hilfe der dazugehörigen Deskriptoren bestimmen. Diese Korrespondenzen erlauben es uns wiederum die Position und Ausrichtung der Anfragekamera zu berechnen. Die Hauptschwierigkeit beim Deskriptorenvergleich liegt dabei in der Größe des betrachteten Problems da unsere Szenenmodelle mehrere Millionen 3D Punkte enthalten während tausende von Features in den Anfragebildern gefunden werden. Als ein Hauptbeitrag dieser Arbeit zeigen wir, dass selbst solche großen Vergleichsprobleme immer noch effizient mittels priorisierten Suchverfahren gelöst werden können. Wir stellen dabei ein einfach umzusetzendes Priorisierungsverfahren vor, welches in der Praxis trotzdem eine nahezu optimale Lösung darstellt. Wir verbinden dabei unsere Priorisierungsstrategie mit einem neuen Ansatz der aktiv nach weiteren Korrespondenzen sucht. Das resultierende Verfahren zur bildbasierten Lokalisierung erreicht dabei die schnellsten Laufzeiten die bisher veröffentlicht wurden während es andere Verfahren in Effektivität erreicht oder sogar übertrifft.

Wir zeigen außerdem, dass die große Stärke dieser Klasse von Verfahren, ihre Fähigkeit qualitativ hochwertige Korrespondenzen zu finden, gleichzeitig deren Anwendbarkeit auf beliebig große Datensätze verhindert. Im letzten Teil der Arbeit beschäftigen wir uns daher mit besser skalierenden Ansätzen und zeigen wie diese Skalierbarkeit mit Effizienz und Effektivität in Einklang gebracht werden kann.



# Acknowledgments

I want to thank my parents for their never-ending love and support. Without them, none of what I have achieved would have been possible.

I thank my advisors, Leif Kobbelt and Bastian Leibe, for providing this opportunity and a nurturing environment in which I could grow both academically and as a person. I am very grateful for all the (technical) discussions and the guidance and help they offered me over all the years. They provided direction when I got lost and they are the reason I fell in love with Computer Vision and Computer Graphics.

I thank all my colleagues at the Computer Graphics and Computer Vision groups, Alex, Alexander, Arne, Darko, David B., David S., Dennis, Dominik, Ellen, Esther, Georgios, H.C., Henrik, Jan, Johannes, Jun, Lars, Lucas, Marcel, Martin, Michael, Mike, Ming, Patrick, Robin, Robert, Sven, Tobias, Volker, and Wolfgang, for the great times that I had over all this years. I will never forget the kart sessions with Darko and Arne, the weekly Friday lunch, and especially the "City Reconstruction" meetings. A special thanks goes to Jan for the excellent technical support and keeping up with me when I again occupied most of the hard disk space or copied massive amounts of data over the internal network. A very special "thank you" goes to Dennis Mitzel for all the fun we had at the conferences we visited together. I am even more thankful towards Tobias Weyand for all our discussions, listening to and improving on my ideas, and helping with collecting datasets.

I want to thank Ole for being the best HiWi ever.

Finally, I want to thank all of my friends, Ann, Bianca, Bert & Manu, Claudia, Dirk & Laura, Elias, Eugen, Henrik, Jan, Jana & Manish, Jan-Thorsten, Mehmet, Melanie, Micha, Michel, Nadine, Oliver, Robert, Sara & Torsten, just to name a few, for their support and patients with me, especially as I too often stayed at work instead of meeting them. They are a major reason I enjoyed my time in Aachen and I will always remember the times we had.

This thesis is dedicated to my brother.



# Contents

<b>1. Introduction</b>	<b>1</b>
1.1. Solving the Image-Based Localization Problem . . . . .	3
1.2. Related Work . . . . .	7
1.3. Contributions & Overview . . . . .	10
<b>2. Foundations</b>	<b>13</b>
2.1. Camera Model . . . . .	13
2.2. Structure-from-Motion . . . . .	16
2.3. Local Features . . . . .	19
2.3.1. RootSIFT . . . . .	21
2.3.2. Compact Binary Descriptor Representations . . . . .	22
2.3.3. Image Retrieval . . . . .	23
<b>I. Feature Matching and Robust Pose Estimation</b>	<b>25</b>
<b>3. Correspondence Search</b>	<b>29</b>
3.1. Approximate Nearest Neighbor Search . . . . .	30
3.1.1. kd-Tree Search . . . . .	32
3.1.2. Hierarchical $k$ -Means Trees . . . . .	32
3.1.3. Other Search Methods . . . . .	34
3.2. Correspondence Search for 3D Models . . . . .	35
3.2.1. Adapting the Ratio Test . . . . .	35
3.2.2. 2D-to-3D vs. 3D-to-2D Matching . . . . .	37
3.3. Quantized Search . . . . .	39
3.4. Discussion . . . . .	41
<b>4. RANSAC-Based Pose Estimation</b>	<b>43</b>
4.1. The N-Point Pose Problem . . . . .	43
4.1.1. Calibrated Cameras . . . . .	44
4.1.2. Uncalibrated Cameras . . . . .	47

4.2. RANSAC . . . . .	49
4.2.1. Introduction to RANSAC . . . . .	51
4.2.2. Spatial Consistent RAnDoM SAmpLe Consensus . . . . .	56
4.2.2.1. Identifying Possible Outliers . . . . .	56
4.2.2.2. The Spatial Consistency Check (SCC) . . . . .	58
4.2.2.3. SCRAMSAC . . . . .	60
4.2.3. Experimental Results . . . . .	62
4.3. Discussion . . . . .	67
 <b>II. Large-Scale Localization using Direct Matching</b>	 <b>71</b>
 <b>5. Fast Direct 2D-to-3D Matching for Image-Based Localization</b>	 <b>75</b>
5.1. 2D-to-3D <i>vs.</i> 3D-to-2D Matching for Localization . . . . .	76
5.1.1. 2D-to-3D Matching . . . . .	77
5.1.2. 3D-to-2D Matching . . . . .	79
5.1.3. Experimental Evaluation . . . . .	80
5.2. Vocabulary-Based Prioritized Search . . . . .	85
5.2.1. A Prioritization Scheme for 2D-to-3D Matching . . . . .	85
5.2.2. Parameter Evaluation . . . . .	90
5.2.3. Comparison With Previous Work . . . . .	96
5.3. Discussion . . . . .	99
 <b>6. Active Correspondence Search for Direct Matching</b>	 <b>101</b>
6.1. Active Correspondence Search . . . . .	102
6.1.1. Prioritization . . . . .	104
6.1.2. Efficient Implementation of Quantized 3D-to-2D Matching . . . . .	106
6.1.3. Computational Complexity . . . . .	108
6.1.4. Discussion of Active Search . . . . .	109
6.2. Visibility Filtering . . . . .	110
6.3. Experimental Evaluation . . . . .	113
6.3.1. Parameter Evaluation . . . . .	116
6.3.2. Faster Linear Search Through Cache Consistency . . . . .	121
6.3.3. Localization Accuracy . . . . .	123
6.3.4. Comparison With State-of-the-Art . . . . .	125
6.4. Discussion . . . . .	129
 <b>III. Scalability of Image-Based Localization Approaches</b>	 <b>131</b>

<b>7. The Scalability of Direct Matching</b>	<b>135</b>
7.1. Limitations of the SIFT Ratio Test . . . . .	136
7.2. Better Descriptor Representations . . . . .	139
7.3. Compact Models for Direct Matching . . . . .	142
7.3.1. Building Compact Models . . . . .	144
7.3.2. Experimental Setup . . . . .	145
7.3.3. Using Compact Models to Reduce Memory Requirements . . . . .	146
7.3.4. Evaluating the Scalability of Compact Models . . . . .	154
7.4. Discussion . . . . .	156
<b>8. Image Retrieval for Scalable Localization</b>	<b>159</b>
8.1. Image Retrieval Revisited . . . . .	161
8.1.1. Image Retrieval for Image-based Localization . . . . .	165
8.1.2. Retrieval <i>vs.</i> Direct Matching . . . . .	167
8.2. Selective Voting . . . . .	168
8.3. Efficient Correspondence Selection . . . . .	171
8.4. Experimental Evaluation . . . . .	173
8.4.1. The Impact of Incorrect Votes . . . . .	174
8.4.2. Correspondence Selection . . . . .	178
8.5. Discussion . . . . .	181
<b>9. Conclusion</b>	<b>183</b>
9.1. Summary & Contributions . . . . .	183
9.2. Future Work . . . . .	186
<b>Bibliography</b>	<b>189</b>