

# *Perceptual Segmentation of Visual Streams by Tracking of Objects and Parts*

---

DISSERTATION

ZUR ERLANGUNG DES MATHEMATISCH-NATURWISSENSCHAFTLICHEN DOKTORGRADES  
"DOCTOR RERUM NATURALIUM"  
DER GEORG-AUGUST-UNIVERSITÄT GÖTTINGEN  
IM PROMOTIONSPROGRAMM  
DER GEORG-AUGUST UNIVERSITY SCHOOL OF SCIENCE (GAUSS)

VORGELEGT VON  
JÉRÉMIE PAPON AUS SUMMIT, NJ, USA



GÖTTINGEN, 2014



# *Perceptual Segmentation of Visual Streams by Tracking of Objects and Parts*

---

DISSERTATION

IN ORDER TO OBTAIN THE DOCTORAL DEGREE IN MATHEMATICS AND NATURAL SCIENCES  
"DOCTOR RERUM NATURALIUM"  
OF THE GEORG-AUGUST-UNIVERSITÄT GÖTTINGEN  
IN THE DOCTORAL PROGRAM OF  
THE GEORG-AUGUST UNIVERSITY SCHOOL OF SCIENCE (GAUSS)

SUBMITTED BY  
JÉRÉMIE PAPON OF SUMMIT, NJ, USA



GEORG-AUGUST-UNIVERSITÄT GÖTTINGEN  
GÖTTINGEN, GERMANY  
OCTOBER 2014



Referentin/Referent: Prof. Dr. Florentin Wörgötter  
Koreferentin/Koreferent: Prof. Dr. Dieter Hogrefe  
Tag der mündlichen Prüfung:

The canonical version of this document is the electronic copy maintained in the Github repository by the author. At this time, it is maintained at:

[https://github.com/jpapon/papon\\_thesis/thesis.pdf](https://github.com/jpapon/papon_thesis/thesis.pdf)

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License. The full terms of the license can be viewed online at:

<http://creativecommons.org/licenses/by-nc/4.0/>

Much of the code created as a result of the research in this thesis is freely available under a BSD license as part of the Point Cloud Library:

<http://www.pointclouds.org/>

The code for the Oculus Vision System (see Appendix A) created as part of this thesis is freely available under GPLv3:

<https://launchpad.net/oculus>

For other usage, contact [jpapon@gmail.com](mailto:jpapon@gmail.com).

# *Perceptual Segmentation of Visual Streams by Tracking of Objects and Parts*

## ABSTRACT

THE ABILITY TO PARSE VISUAL STREAMS into semantically meaningful entities is an essential element of intelligent systems. This process - known as segmentation - is a necessary precursor to high-level behavior which uses vision, such as identification of objects, scene understanding, and task planning. Tracking these segmented entities over time further enriches this knowledge by extending it to the action domain. This work proposes to establish a closed loop between *Video Object Segmentation (VOS)* and *Multi-Target Tracking (MTT)* to parse streaming visual data. We demonstrate the strengths of this approach, and show how such a framework can be used to distill basic semantic understanding of complex actions in real-time, without the need for a-priori object knowledge. Importantly, this framework is highly robust to occlusions, fast movements, and deforming objects.

This thesis has four key contributions, each of which lead towards fast and robust video segmentation through tracking. First, we present *Video Segmentation by Relaxation of Tracked Masks (VSRTM)*, which serves as a proof of concept, demonstrating the feasibility of Dynamic Semantic Segment Tracking (DSST) in 2D video. This method serves as a demonstration of the viability of a feedback loop between VOS and MTT. This is accomplished using a sequential Bayesian technique to generate predictions which are used to seed a segmentation kernel, the results of which are used to update tracked models.

The second contribution consists of a 3D voxel clustering technique, *Voxel Cloud Connectivity Segmentation (VCCS)*, which makes use of a novel adjacency octree structure to efficiently cluster 3D point cloud data, and provide a graph lattice for the otherwise unstructured points. These clusters of voxels, or *supervoxels*, and their adjacency graph are used to maintain a world model which serves as an internal buffer for observations for trackers. Importantly, this world model uses ray-tracing to ensure that it does not delete occluded voxels as new frames of data arrive.

The third contribution is a novel spatially stratified sampling technique for evaluating the likelihood function in particle filters. In particular, we show that in the case where the measurement function uses spatial correspondence, we can greatly reduce computational cost by exploiting spatial structure to avoid redundant computations. We present results which quantitatively show that the technique permits equivalent, and in some cases, greater accuracy, as a reference point cloud particle filter at significantly faster run-times. We also compare to a GPU implementation, and show that we can exceed their performance on the CPU. In addition, we present results on a multi-target tracking application, demonstrating that the increases in efficiency permit online 6DoF multi-target tracking on standard hardware.

Our final contribution is *Predictive Association of Supervoxels* (*PAS*), which implements a closed loop between segmentation and tracking by minimizing a global energy function which scores supervoxel associations. The energy function is efficiently computed using the adjacency octree, with candidate associations provided by the 3D correspondence based particle filters. The association found determines a fully segmented point cloud, and is used to update the tracker models (as in VSRTM). This allows for the segmentation of temporally consistent supervoxels, avoiding the need to pre-define object models for segmentation.

Each of these contributions has been implemented in live systems and run in an online streaming manner. We have performed quantitative evaluation on existing benchmarks to demonstrate state-of-the-art tracking and segmentation performance. In the 2D case, we compare against an existing tracking benchmark, and show that we can match their tracking performance, while in the 3D case we use a benchmark to show that we can outperform a GPU implementation. Finally, we give qualitative results in a robotic teaching application, and show that the system is able to parse real data and to distill semantic understanding from video.

# Contents

1	INTRODUCTION	1
1.1	Problem Definition and Motivation . . . . .	2
1.1.1	The Image Segmentation Problem . . . . .	2
1.1.2	The Tracking Problem . . . . .	3
1.1.3	Video Object Segmentation - Segmentation In Sequential Frames .	5
1.2	State of the Art . . . . .	6
1.2.1	Segmentation and Superpixels . . . . .	6
1.2.2	Multi-Target Visual Tracking . . . . .	6
1.2.3	Video Object Segmentation . . . . .	7
1.3	Outline and Contributions . . . . .	8
2	VIDEO SEGMENTATION BY RELAXATION OF TRACKED MASKS	11
2.1	Overview of the Algorithm . . . . .	12
2.2	Tracking Object Masks . . . . .	13
2.2.1	Sequential Bayesian Estimation . . . . .	14
	Dynamic Model . . . . .	14
	Measurement Model . . . . .	15
2.2.2	Parallel Particle Filters . . . . .	15
2.2.3	Particle Birth, Repopulation, & Decay. . . . .	16
2.3	Extracting a Dense Image Labeling . . . . .	17
2.3.1	Object Pixel Likelihood Maps. . . . .	17
2.3.2	Label Association Likelihood Map. . . . .	17
2.4	Occlusion Handling. . . . .	18
2.5	Segmentation using Superparamagnetic Clustering . . . . .	18
2.6	Experimental Results . . . . .	20
2.7	Discussion . . . . .	21
3	PATCH-BASED PERCEPTUAL WORLD MODEL	25

3.1	Pre-processing of Point Cloud Data . . . . .	26
3.1.1	Voxelization . . . . .	26
3.1.2	Octree Adjacency Graph . . . . .	26
3.2	Geometrically Constrained Supervoxels . . . . .	27
3.2.1	Spatial Cluster Seeding . . . . .	28
3.2.2	Cluster Features and Distance . . . . .	29
3.2.3	Flow Constrained Region Growing . . . . .	30
3.3	Sequential Update of Perceptual Model . . . . .	31
3.4	Depth Dependent Voxel Grid . . . . .	33
3.5	Locally Convex Connected Patches . . . . .	34
3.6	Experimental Results . . . . .	37
3.6.1	Datasets . . . . .	37
	Object Segmentation Database (OSD) . . . . .	37
	NYU Indoor Dataset (NYU) . . . . .	37
	Returning to the Projected Plane . . . . .	38
3.6.2	Supervoxels . . . . .	40
	Object Boundary Adherence . . . . .	40
	Time Performance . . . . .	41
3.6.3	Locally Convex Connected Patches . . . . .	42
3.7	Discussion . . . . .	43
4	<b>MODEL-BASED POINT CLOUD TRACKING</b>	<b>45</b>
4.1	Particle Filters in 3D . . . . .	46
4.1.1	Model Representation . . . . .	46
4.1.2	Dynamic Model . . . . .	47
4.1.3	Measurement Model . . . . .	48
4.2	Stratified Correspondence Sampling . . . . .	50
4.3	Experimental Results . . . . .	51
4.3.1	Results on Synthetic Sequences . . . . .	52
4.3.2	Results on Real Sequences . . . . .	57
4.4	Discussion . . . . .	59
5	<b>TRACKING BASED POINT CLOUD VIDEO SEGMENTATION</b>	<b>61</b>
5.1	Tracked Model Representation . . . . .	62
5.2	Bank of Parallel Particle Filters . . . . .	63
5.3	Association by Joint Label Optimization . . . . .	63
5.4	Alignment and Update of Models . . . . .	65
5.5	Experimental Results . . . . .	66

5.5.1	Imitation of Trajectories for Robot Manipulation . . . . .	66
5.5.2	Semantic Summaries of Actions . . . . .	68
5.6	Discussion . . . . .	68
<b>6</b>	<b>CONCLUSIONS</b>	<b>71</b>
6.1	Summary of Contributions . . . . .	71
6.2	Shortcomings of VOS Benchmarks . . . . .	73
6.3	Limitations and Direction of Future Work . . . . .	73
<b>REFERENCES</b>		<b>81</b>
<b>APPENDICES</b>		<b>83</b>
<b>A</b>	<b>THE OCULUS VISION SYSTEM</b>	<b>85</b>
A.1	Motivation . . . . .	85
A.2	System Architecture . . . . .	86
A.2.1	Execution Flow . . . . .	86
A.2.2	Plugin Development and Interaction . . . . .	87
A.2.3	Visualization . . . . .	89
A.3	Memory Architecture . . . . .	89
A.3.1	Global Buffer . . . . .	89
A.3.2	GPU Memory Handling . . . . .	91
A.4	Demonstration System . . . . .	92
A.4.1	Image Acquisition . . . . .	92
A.4.2	Disparity and Optical Flow . . . . .	93
A.4.3	Segmentation and Tracking . . . . .	93
A.4.4	Semantic Graphs . . . . .	94
A.5	Results and Discussion . . . . .	95
A.6	Conclusion . . . . .	96
<b>B</b>	<b>SEQUENTIAL BAYESIAN ESTIMATION</b>	<b>97</b>
B.1	Particle Filters . . . . .	98
B.1.1	Resampling . . . . .	98

# List of Figures

## 1 Introduction

1.1.1 Example of Segmentation and Ground Truth . . . . .	3
1.1.2 Technical Difficulties of Segmentation . . . . .	4
1.1.3 Hidden Markov Model . . . . .	4
1.1.4 Example of Visual Tracking . . . . .	5
1.1.5 Example of Video Object Segmentation . . . . .	6

## 2 Video Segmentation by Relaxation of Tracked Masks

2.1.1 Overview of Algorithm . . . . .	13
2.5.1 Relaxation Convergence . . . . .	19
2.6.1 Tracked output from lemming sequence . . . . .	22
2.6.2 Results of Cranfield Sequence . . . . .	23

## 3 Patch-based Perceptual World Model

3.1.1 Example of Voxelization . . . . .	26
3.1.2 Octree Voxelization . . . . .	27
3.1.3 Adjacency in a 3d Grid . . . . .	27
3.2.1 Seeding Parameters . . . . .	29
3.2.2 Seeding Size . . . . .	29
3.2.3 Voxel Search Order . . . . .	31
3.3.1 Voxel Visibility . . . . .	32
3.3.2 Voxel Permanence . . . . .	33
3.4.1 Depth Adaptive Transform . . . . .	34
3.5.1 Flow Diagram of LCCP . . . . .	36
3.6.1 NYU Dataset Examples . . . . .	38
3.6.2 2D Hole Filling . . . . .	39
3.6.3 Superpixel Comparison . . . . .	39
3.6.4 Boundary Recall & Undersegmentation Error . . . . .	40
3.6.5 Segmentation Speed . . . . .	41
3.6.6 OSD Dataset Examples . . . . .	42

<b>4 Model-Based Point Cloud Tracking</b>	
4.1.1 Example of data from “Tide” sequence. . . . .	46
4.2.1 Stratified Correspondence Matching . . . . .	50
4.2.2 Tracking on the artificial “Kinect Box” sequence. . . . .	51
4.3.1 Tracking on the artificial “Tide” sequence. . . . .	53
4.3.2 Tracked vs Ground Truth - Kinect Box . . . . .	54
4.3.3 Results on the Kinect Box artificial sequence. . . . .	54
4.3.4 Tracked vs Ground Truth - Milk . . . . .	55
4.3.5 Results on the Milk artificial sequence. . . . .	55
4.3.6 Tracked vs Ground Truth - Orange Juice . . . . .	56
4.3.7 Results on the Orange Juice artificial sequence. . . . .	56
4.3.8 Tracked vs Ground Truth - Tide . . . . .	57
4.3.9 Results on the Tide artificial sequence. . . . .	57
4.3.10 Human demonstration of assembly of the Cranfield Scenario. . . . .	58
4.3.11 Snapshots from Virtual Reality Benchmark Run . . . . .	59
<b>5 Tracking Based Point Cloud Video Segmentation</b>	
5.1.1 Algorithm Overview . . . . .	62
5.1.2 The Aperture Problem . . . . .	63
5.3.1 Supervoxel Association . . . . .	64
5.5.1 Cranfield Tracking Results . . . . .	66
5.5.2 Trajectory Imitation . . . . .	67
5.5.3 Cranfield Key Frames . . . . .	68
<b>A The Oculus Vision System</b>	
A.2.1 Overview of the system architecture . . . . .	88
A.3.1 Comparison of Buffering Schemes . . . . .	90
A.3.2 Feedback using a Global Buffer . . . . .	91
A.3.3 Streaming and Concurrent Kernels . . . . .	92
A.4.1 Timing results for demonstration system . . . . .	94
A.5.1 Performance Effect of Visualization . . . . .	96

# List of Tables

2.6.1 PROST dataset benchmark results . . . . .	20
3.6.1 Segmentation Results on OSD Dataset . . . . .	42
3.6.2 Comparison of NYU Dataset Results . . . . .	43

# List of Acronyms

**AI** Artificial Intelligence.

**DDVG** Depth Dependent Voxel Grid.

**DSST** Dynamic Semantic Segment Tracking.

**ECC** Extended Convexity Criterion.

**LCCP** Locally Convex Connected Patches.

**MHVS** Multiple hypothesis video segmentation.

**MSVS** Mean-shift video segmentation.

**MTT** Multi-Target Tracking.

**MTVT** Multi-target visual tracking.

**PAS** Predictive Association of Supervoxels.

**PDF** Probability Distribution Function.

**PVA** Propagation, validation, and aggregation.

**SBF** Sequential Bayesian Filtering.

**VCCS** Voxel Cloud Connectivity Segmentation.

**VOS** Video Object Segmentation.

**VSRTM** Video Segmentation by Relaxation of Tracked Masks.



# List of Related Publications

**Papon, J.**; Wörgötter, F., “Spatially Stratified Correspondence Sampling for Real-Time Point Cloud Tracking,” *Applications of Computer Vision (WACV), 2015 IEEE International Conference on*, Jan. 2015.

**Papon, J.**; Kulvicius, T.; Aksoy, E.; Wörgötter, F., “Point Cloud Video Object Segmentation using a Persistent Supervoxel World-Model,” *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, Nov. 2013.

**Papon, J.**; Abramov, A.; Schoeler, M.; Wörgötter, F., “Voxel Cloud Connectivity Segmentation - Supervoxels for Point Clouds,” *Computer Vision and Pattern Recognition (CVPR) 2013*, June 2013.

**Papon, J.**; Abramov, A.; Wörgötter, F., “Occlusion Handling in Video Segmentation via Predictive Feedback,” *European Conference on Computer Vision (ECCV) 2012. Workshops and Demonstrations, Lecture Notes in Computer Science Volume 7585, 2012*, pp 233-242.

**Papon, J.**; Abramov, A.; Aksoy, E.; Wörgötter, F., “A modular system architecture for online parallel vision pipelines,” *Applications of Computer Vision (WACV) 2012*, pp.361-368, Jan. 2012.

Stein, S.; Schoeler, M.; **Papon, J.**; Wörgötter, F., “Object Partitioning using Local Convexity,” *Computer Vision and Pattern Recognition (CVPR) 2014*, June 2014.

Stein, S.; Wörgötter, F.; Schoeler, M.; **Papon, J.**; Kulvicius, T., “Convexity Based Object Partitioning For Robot Applications,” *Robotics and Automation (ICRA), 2014 IEEE/RSJ International Conference on*, June 2014.

Schlette, C.; Buch, A.; Aksoy, E.; Steil, T.; **Papon, J.**; Savarimuthu, T.R.; Wörgötter, F.; Krüger, N.; Roßmann, J., “A new benchmark for pose estimation with ground truth

from virtual reality,” *Production Engineering*, May 2014.

Aein, M.J.; Aksoy, E.; Tamosuinaite, M.; **Papon, J.**; Ude, A.; Wörgötter, F., “**Toward a library of manipulation actions based on semantic object-action relations**,” *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, Nov. 2013.

Abramov, A.; Pauwels, K.; **Papon, J.**; Wörgötter, F.; Dellen, B., “**Depth-supported real-time video segmentation with the Kinect**,” *Applications of Computer Vision (WACV) 2012*, Jan. 2012.

The research leading to this thesis was supported with funding from the European Community’s Seventh Framework Programme FP7/2007-2013 (Specific Programme Cooperation, Theme 3, Information and Communication Technologies) under grant agreement no. 270273, Xperience and grant agreement no. 269959, IntellAct.

“Le seul véritable voyage, ce ne serait pas d’aller vers de nouveaux paysages, mais d’avoir d’autres yeux, de voir l’univers avec les yeux d’un autre, de cent autres, de voir les cent univers que chacun d’eux voit, que chacun d’eux est.”

## Acknowledgments

*The only true voyage, would not be a journey towards new landscapes, but to possess new eyes, to see the universe through the eyes of another, of a hundred others, to behold the hundred universes that each of them sees, that each of them is.*

Marcel Proust

I WOULD NEVER HAVE MADE IT THROUGH to the end of this thesis without the support of my friends, family, and colleagues. I’d like to thank my supervisors Prof. Dr. Florentin Wörgötter and Prof. Dr. Dieter Hogrefe for their guidance, aid, and many fruitful discussions which helped bring about this work. Special thanks go out to our close-knit (talkative) vision group: Dr. Eren Erdal Aksoy, Simon Reich, Simon Stein, and Markus Schöler. I’d also like to thank our robot-men Dr. Tomas Kulvicius and Mohamad Javad Aein for all their hard work making it work (usually) in the real-world. I also would like to thank all of my friends and colleagues up at the University of Southern Denmark, as well as all the IntellAct partners with whom I had the good fortune of passing so much time in Odense. A heart-felt thanks to all of the mem-

bers of Florentin’s group, which has become my extended family in Germany. It’s been a pleasure working and living with all of you: Mohamad Javad Aein, Dr. Alejandro Agostini, Martin Biehl, Jan-Matthias Braun, Sakyasingha Dasgupta, Michael Fauth, Dennis Goldschmidt, Dr. Yinyun Li, Timo Nachstedt, Dr. Poramate Manoonpong, Dr. Minija Tamosiunaite, Dr. Christian Tetzlaff, and Xiaofeng Xiong. An especially big thanks to Ursula Hahn-Wörgötter for putting up with me and being such a big help in figuring out Germany. Last (but far from

least) I want to thank my family. I could have never made it here without the unwavering support of my father, Jean-Marc, and my mother, Marian. I thank them especially for hosting me every Summer in Veyssou and on the boat - time away from work that proved invaluable. I know at times I was a handful to deal with (perhaps a few handfuls), and your constant love and support were instrumental in seeing me through. Of course I would also like to thank my loving sister Camille, who, even half a world away, hasn’t forgotten her little brother.

Thank you all, so very, very much!

JÉRÉMIE PAPON

GÖTTINGEN, 2014.

*We are so familiar with seeing, that it takes a leap of imagination to realize that there are problems to be solved. But consider it. We are given tiny distorted upside-down images in the eyes, and we see separate solid objects in surrounding space. From the patterns of stimulation on the retina we perceive the world of objects and this is nothing short of a miracle.*

Richard L. Gregory, *Eye and Brain*, 1966.

# 1

## Introduction

The human visual cortex is able to process a bewilderingly large amount of data with ease. From messy signals emitted by the 100 million rods and cones in a typical retina, it can assemble an ordered world containing structure, meaningful parts, and distinct objects [45]. Furthermore, it possesses an understanding of coherent motion, allowing it to keep track of and intuitively predict object trajectories. These two abilities, the segmentation of the world into objects, and the tracking of objects to maintain their identities, serve as key components in the bootstrapping of higher level knowledge. Indeed, it has been shown that our earliest and most fundamental understanding of the world is topological in nature, dealing with concepts that can be described through segmentation and tracking - proximity, order, separation and enclosure[66].

In fact, these concepts are so fundamental to human understanding of the world that we find it profoundly difficult to precisely define what an object actually is. Yet in spite of the difficulty in formalizing the concept, we can divide complex moving scenes into distinct objects, even hierarchies of parts, with little effort. In this work we argue that the concepts of tracking and segmentation are inexorably linked; that visual tracking plays an essential role in creating the objects we observe, and that the organization of observations into structured objects is critical for robust tracking. We propose that without the ability to track motions in a coherent way, the notion of distinct objects is, ultimately, a meaningless one. Furthermore, we suggest that this link between tracking and object segmentation is one of the key elements that enable learning from visual input, and through this, the bootstrapping of cognition itself.

## 1.1 PROBLEM DEFINITION AND MOTIVATION

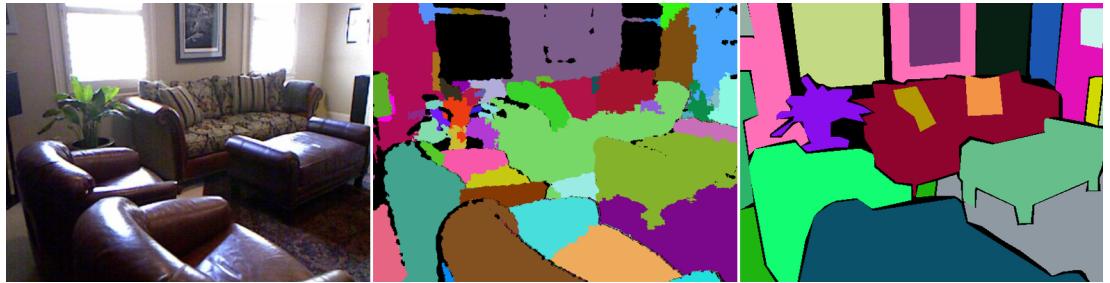
As with humans, in order for intelligent agents to be truly autonomous, they must be able to learn the principles of visual understanding from their own unsupervised observations. At its most basic level, an agent must be able to parse observations, to break them down into meaningful entities upon which higher level knowledge can be built. In other words, segmentation of observations is a precursor to high-level behaviors, such as identification of objects, scene understanding, and task planning. Tracking segmented entities over time is an integral parts of this, as it further enriches this knowledge by extending it to the action domain. Combining these two tasks - segmentation and tracking - would allow fully unsupervised parsing of streaming visual data. This has the potential to greatly increase the flexibility of autonomous robotic systems by allowing them to learn from observations without the constraints of pre-defined object and domain knowledge.

In this work, we propose to develop an unconstrained video segmentation algorithm that is able to track low level patches. This permits the segmentation of objects and their parts naturally, without the need to define what an object actually is. Rather than train classifiers to recognize pre-defined objects, we can have an agent observe or interact with a scene and learn the concept of an object through movement and interactions between observed patches. This Chapter introduces the general concepts that will be expounded upon throughout the work by first discussing the three underlying tasks; Image Segmentation, Multi-Target Tracking (MTT), and Video Object Segmentation (VOS). With each of these tasks, we will discuss what exactly our goals are, and what challenges are faced in achieving them. Next, we survey the state of the art in each of these fields, highlighting the methods and important papers upon which we base this work. Finally, we outline each of the Chapters of this work, and enumerate the specific contributions of our research.

### 1.1.1 THE IMAGE SEGMENTATION PROBLEM

Image segmentation aims to divide the set of pixels in an image into a number of distinct subsets, where each subset represents some semantically meaningful entity (e.g., an object - see Figure 1.1.1). This is a (infamously) deceptively tricky business, primarily because it is something that humans are able to do intuitively. This ease with which humans can segment visual scenes is highly deceptive; Marvin Minsky, one of the pioneers of Artificial Intelligence (AI), famously assigned one of his students “computer vision” as a summer undergraduate project in 1966. Nearly half a century later, despite the extensive effort to solve it, image segmentation, the first step on the long road to complete “computer vision”, remains an unsolved problem. In fact, this phenomenon - of tasks that are simple for humans being incredibly demanding computationally - even has its own name; *Moravec’s Paradox*. As stated by Pinker [67]:

“The main lesson of thirty-five years of AI research is that the hard problems are



**Figure 1.1.1:** Example of Segmentation and one interpretation of Ground Truth. From left to right we have an image, a segmentation from a computer vision algorithm, and a human-annotated ground truth labeling. Here labels are represented by different colors, a convention we shall use throughout the rest of this work.

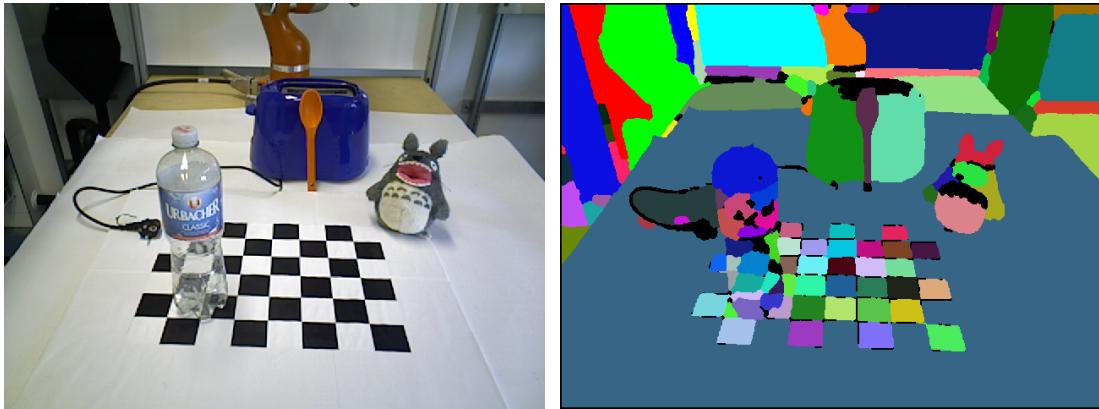
easy and the easy problems are hard. The mental abilities of a four-year-old that we take for granted – recognizing a face, lifting a pencil, walking across a room, answering a question – in fact solve some of the hardest engineering problems ever conceived. (p. 190)”

The reason for this “hardness” of an “easy” problem like image segmentation is two-fold: firstly, there are many technical and computationally-demanding challenges associated with properly dividing an image into separate objects. Among these, shadows, occlusions, reflections, imaging noise and so forth can all greatly affect the results of image segmentation. Consider, for instance, a partial occlusion as in Figure 1.1.2. A human can easily identify that the parts on either side of the occluding object belong to same object. This is accomplished using what we shall refer to as *high-level* knowledge throughout this work - in this case, knowledge of the complete nature of an object.

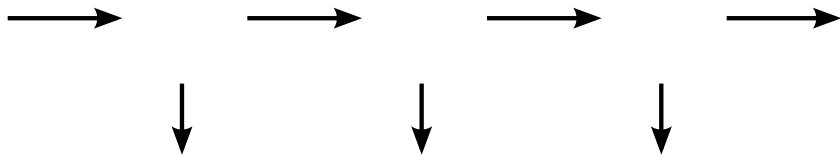
This leads us to the second challenge in image segmentation, which is that, generally speaking, there is no “correct” solution to the problem. A perfect labeling for one application might be useless in another. This is even more of a problem when we are discussing segmentation separate from any application, as is the case with standard image segmentation benchmarks (which are used to quantify algorithm performance). These benchmarks use ground-truth image labels (manually created by humans) to score the output of different algorithms. Unfortunately, the correctness of different labellings is highly subjective, and hand-drawn labels from people can differ radically.

### 1.1.2 THE TRACKING PROBLEM

Tracking entities over time is a critical element in a wide variety of computer vision applications such as visual surveillance, action recognition, and robotic imitation learning. In most of these, visual tracking serves as the precursor to further high-level inference, as without it, one is unable to correctly interpret time-variant systems. One can formalize the tracking problem



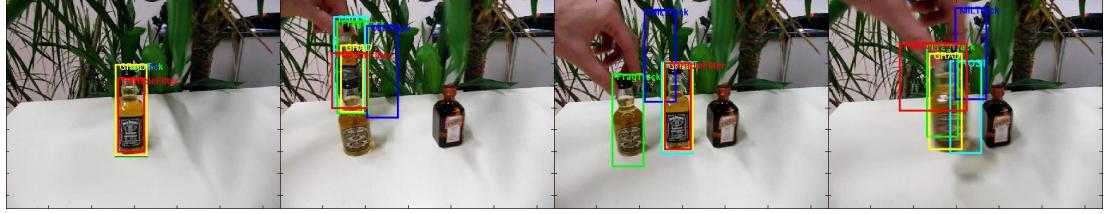
**Figure 1.1.2:** Technical Difficulties of Segmentation. Here we see some of the myriad of technical difficulties present in color-based segmentation, such as transparent objects (the water bottle), partial occlusions (the toaster), objects with strong color differences (the little monster), and similarities in color (the bottle cap to the table).



**Figure 1.1.3:** The Hidden Markov Model is a classical way to represent the track of an object over time. The object states  $x(0 \dots t)$  (shown here in blue) are hidden variables which influence observations  $y(0 \dots t)$  through conditional dependencies (shown as arrows). An important property of the Markov Model is that state at time  $t$  is dependent only on state at time  $t - 1$ .

as estimation of the time-varying hidden state (e.g., position, velocity) of an object  $x(t)$  using noisy observations  $y(t)$ . For simplicity, one generally assumes the state evolution to be a *Markov Process* (see Figure 1.1.3), that is, a stochastic process which is conditionally independent of the rest of its history given its previous state.

Multi-target visual tracking (MTVT) extends these concepts to multiple targets, adding additional complexity due to the need to both estimate the number of tracked targets as well as associate observations with appropriate targets. This is the primary challenge of MTVT - the data association problem - deciding which tracked target a particular observation belongs to. Confounding this is the additional null possibility, where an observation belongs to none of the tracked targets. Some additional difficulties present in MTVT are related to those of image segmentation, simply extended into the temporal domain. In particular, interacting and occluded targets are especially challenging.



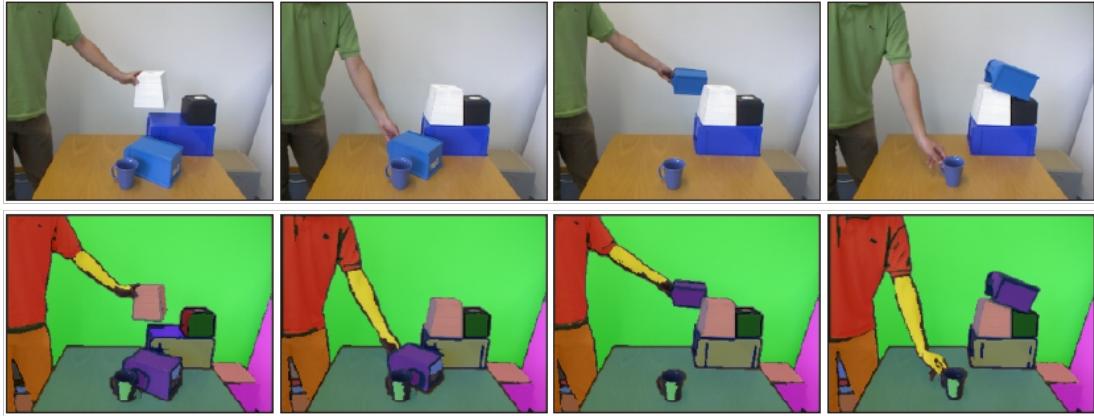
**Figure 1.1.4:** Example of Visual Tracking - from [61]. This shows outputs from various trackers in a standard video tracking benchmark [76]. The output of each tracker is shown as a colored rectangle. Some of the difficulties of tracking can be seen- in particular complex backgrounds, motion blur, partial occlusions (second frame from left) and even full occlusions (right-most frame).

### 1.1.3 VIDEO OBJECT SEGMENTATION - SEGMENTATION IN SEQUENTIAL FRAMES

VOS attempts to cluster pixels of video frames into segments which are both spatially and temporally coherent. While related to MTVT, VOS goes a step beyond localizing tracked objects, in that it makes an association decision for each observed pixel; in addition to estimating overall state, it must re-estimate spatial extent every frame. Additionally, VOS has the additional consideration that target appearance models are unknown a-priori, and are subject to arbitrary changes over time.

The standard interpretation of VOS is that of adding an additional dimension to image segmentation; that is, one stacks all the image frames on top of each other, and performs a “volumetric” segmentation. In this work we shall use a different interpretation for VOS; that of tracking multiple time-evolving and interacting objects projected onto the image plane of our sensor. While the standard interpretation has the advantage of allowing the straightforward extension of 2D segmentation techniques, it suffers greatly from its inability to handle occlusions in a meaningful way. This is easily observed when one considers that occlusions will result in “disconnection” within the 3D stack, violating the core assumption that segments of interest form contiguous volumes. In contrast, tracking techniques are able to handle occlusions gracefully.

One interesting aspect of video segmentation is that it has the potential to be more accurate than single image segmentation, as it can take advantage of the temporal coherence of objects to infer information about the objects in a scene. Unfortunately, the addition of the temporal domain brings along new challenges as well; for instance that pixels which should be grouped across time may not be continuously visible, as in the case of partial or full occlusions. Additionally, the added dimension increases the computational complexity of the problem, making accurate segmentation a costly procedure. Temporal information also increases the exposure of the algorithm to noise, as each image frame is a separate noisy measurement. This adds a large amount of uncertainty to the problem, since measured values (i.e., of color) for an object can show significant variation over time.



**Figure 1.1.5:** Example of Video Object Segmentation - from [3]. This shows the goal of VOS - to extract a dense labeling (labels here are shown as distinct colors) for every frame, maintaining temporal consistency of objects. For many applications it is of vital importance to make the labeling consistent from frame to frame, that is, to maintain object identities.

## 1.2 STATE OF THE ART

### 1.2.1 SEGMENTATION AND SUPERPIXELS

Segmentation of scenes into objects remains one of the most challenging topics of computer vision despite decades of research. To address this, recent methods often use hierarchies which create a rank order that build bottom-up from small localized superpixels to large-scale regions [7, 11, 71]. As an alternative, researchers have also pursued strictly top-down approaches. Such methods began with coarse segmentations using multiscale sliding window detectors [87], later progressing to finer grained segmentations and detections based on object parts [18, 31]. These two avenues of research led naturally to methods which *combine* bottom-up hierarchy building with top-down object- and part-detectors [12, 37, 79]. While these approaches have yielded quite good results even on complex, varied data sets, they have lost much of the generality of learning-free approaches. In general the most powerful methods to date use trained classifiers for segmentation [37, 79]. This means they cannot be applied to arbitrary unknown scenes without being retrained, requiring the acquisition of a new data-set tailored to each test environment and a-priori models specialized to this testing data.

### 1.2.2 MULTI-TARGET VISUAL TRACKING

MTVT is a well-established field, which goes back over thirty years [32]. In this work we use Sequential Bayesian Estimation to track targets, in particular a Monte Carlo method known as Particle Filtering. This approach was first introduced to the vision community by Isard and Blake [44] and has been the subject of much subsequent research extending it [40, 86, 88].

There are two standard approaches that have been used to extend the Particle Filter to multiple targets. The first represents all targets jointly in a single particle filter by assigning individual particles to particular labels [85]. This means that, for a given total number of particles, there

will be fewer for each individual target - resulting in reduced accuracy. The second approach is to add additional dimensions to the state space for each additional target [77]. Unfortunately, this approach quickly increases the dimensionality of the state space, which also results in a need for a very high number of particles for the filter to remain accurate.

In both of the above approaches, the computational complexity increases exponentially as targets are added (for constant level of accuracy). As a consequence of this, it is beneficial to use a separate particle filter for each target. One way of doing this is to add factors to the observation and/or process models of the filters which explicitly model occlusions and interactions between targets [46, 52]. Alternatively, one can use a discrete processing step to resolve the association of target detections [48].

A different approach which has generated much interest is to use the output of detectors as the basis for tracking. Known as *tracking-by-detection*, these methods typically use simple particle filters to maintain tracks [20, 24], and shift the focus of the problem onto the data association step, wherein detections are assigned to targets. While there are several classical approaches for solving this association problem from Sonar and Radar research [33, 70], a greedy approach is typically sufficient given a good association scoring function [20, 90].

### 1.2.3 VIDEO OBJECT SEGMENTATION

There are many existing VOS methods, which can be classified based on three parameters; whether they are on- or off-line, whether they are dense or sparse, and whether or not they are supervised. We can reduce the comparison-space of related work by comparing only with algorithms which have the same three parameters as this work - on-line processing (the algorithm may only use past data), dense segmentation (every pixel is assigned to a spatio-temporal cluster), and unsupervised operation. Four state-of-the-art segmentation algorithms meet these requirements: Mean-shift video segmentation (MSVS) [64], Multiple hypothesis video segmentation (MHVS) from superpixel flows [83], Propagation, validation, and aggregation (PVA) of a preceding graph [55], and Matching images under unstable segmentations [39]. Of these methods, none are able to handle full occlusions; in fact only MHVS considers occlusions, and it is only able to handle partial occlusions for a few frames, and does not consider full occlusions. Even state of the art off-line methods such as that of Brendel and Todorovic [21] only handle partial occlusions, claiming that “complete occlusions ... require higher-level reasoning”.

In [58] Papadakis and Bugeau use a dynamical model to guide successive segmentations, along with an energy function minimized using graph cuts to solve the label association problem. They formally model visible and occluded regions of tracked objects, tracking them as distinct parts. While they do consider occlusions, they do not maintain a world model, and as such their methodology must fail under complete occlusions. Additionally, they formally model visible and occluded parts of the tracked objects, and so the method does not scale well with an increasing number of objects, and thus is better suited to extracting the silhouettes of

a few objects than performing a full segmentation. Other methods, such as [1], are severely limited in that they require pre-computed models which are calibrated to a ground plane in order to resolve occlusions. Recent work in MTVT [57] successfully tracks multiple objects using a segmentation and association approach and adaptive 3D appearance models, but is limited by the need to align model point clouds to the observed data every frame, as well as the need for a ground plane. This precludes it from handling occlusions, as once a target is no longer observed, its track must be terminated.

### 1.3 OUTLINE AND CONTRIBUTIONS

This work is organized as follows: First, in Chapter 2 we present a hybrid VOS / MTT technique for 2D data. We describe the segmentation algorithm used, how we track segments, how we combine tracked results into a video segmentation and finally present results on a tracking benchmark. In Chapter 3 we present the concept of a persistent 3D voxel world model. We begin by briefly introducing some core concepts of acquisition and representation of 3D point cloud data, then present VCCS, a method for extracting a graph of 3D voxel patches from point cloud data. We then discuss how to add point clouds sequentially to the model in a way that allows voxels to persist through occlusions. Finally, we present quantitative and qualitative results of VCCS and Locally Convex Connected Patches (LCCP), a segmentation method which uses VCCS. In Chapter 4 we describe a method for using particle filters to track multiple rigid objects in point cloud video data and present results of tracking performance on both real and artificial data. Additionally, we present a stratified sampling approach which greatly reduces the computational complexity of tracking. In Chapter 5 we combine the methods described in prior Chapters into a system which can produce full video segmentation of point cloud videos. We show that the system is highly robust to occlusions and noisy data, and present results on the application of semantic understanding and imitation of human actions. Finally, in Chapter 6 we discuss the findings and experimental results of this work, possible future work, and conclude.

Each of the Chapters in this thesis contain novel contributions to the field, briefly described below.

- **Chapter 2** contains a 2D segmentation through relaxation technique published in [61]. This work demonstrated the concept of extracting video segmentation from tracks, and the idea of connecting segmentation and tracking in a closed feedback loop.
- **Chapter 3** contains the Supervoxel clustering method VCCS, as well as the scheme for maintaining voxels in an octree through occlusions published in [62]. Supervoxels serve as the basis for much ongoing work, as they provide a graph structure for otherwise unordered pointcloud data.
- **Chapter 4** accelerates 3D correspondence particle filter tracking through a stratified

sampling of the model-space published in [59]. This technique greatly reduces the computational complexity of pointcloud tracking by taking advantage of the spatial structure of points.

- **Chapter 5** has the techniques used to generate full segmentations based upon the results from multiple independent trackers [63].
- **Appendix A** presents the Oculus Vision System [60], an open-source computer vision system created over the course of the research for this thesis.

Additionally, the methods presented in this work have all been published as open-source and are publicly available, either as part of Oculus<sup>1</sup> or the Point Cloud Library (PCL)<sup>2</sup>.

---

<sup>1</sup><https://launchpad.net/oculus/>

<sup>2</sup><http://www.pointclouds.org/>



*The outcome of any serious research can only be to make two questions grow where only one grew before.*

Thorstein Veblen

# 2

## Video Segmentation by Relaxation of Tracked Masks

**I**N THE BEGINNING, 3D data, especially video data, was not readily available. As such, researchers were forced to make due with strictly 2D video, which is inherently ambiguous in many situations. In particular, partial and full occlusions are particularly vexing problems in 2D video - not least because understanding of 2D video is so easy for humans, yet so difficult to interpret algorithmically. Indeed, knowledge of object permanence, that is, the understanding of how to correctly interpret occlusions, is something that humans acquire very early on in their lives [45], but has yet to be successfully implemented in a fully automated VOS system. Even after decades of research, state-of-the-art methods still have trouble correctly resolving partial occlusions, and typically fail completely after even the briefest of complete occlusions.

In this Chapter, we shall present our attempts towards resolving the object permanence problem with 2D data, as well as advance color-based VOS in general. In particular, we seek to overcome two of the main drawbacks of the color-based video segmentation method developed by Abramov et al. [2] (and indeed, of color-based VOS in general). The first of these is the correct tracking of objects through partial and full occlusions, which we proposed to solve using a layering of deformable object masks that are allowed to interact and compete for “ownership” of pixels. The second is to allow for object identities to be maintained through sudden and/or fast movements - something that was not possible due to the core assumptions

of the algorithm. To correct for this, we tracked the masks with a set of particle filters, a class of Bayesian predictive filters which are well known for their ability to handle difficult trajectories [40, 86, 88].

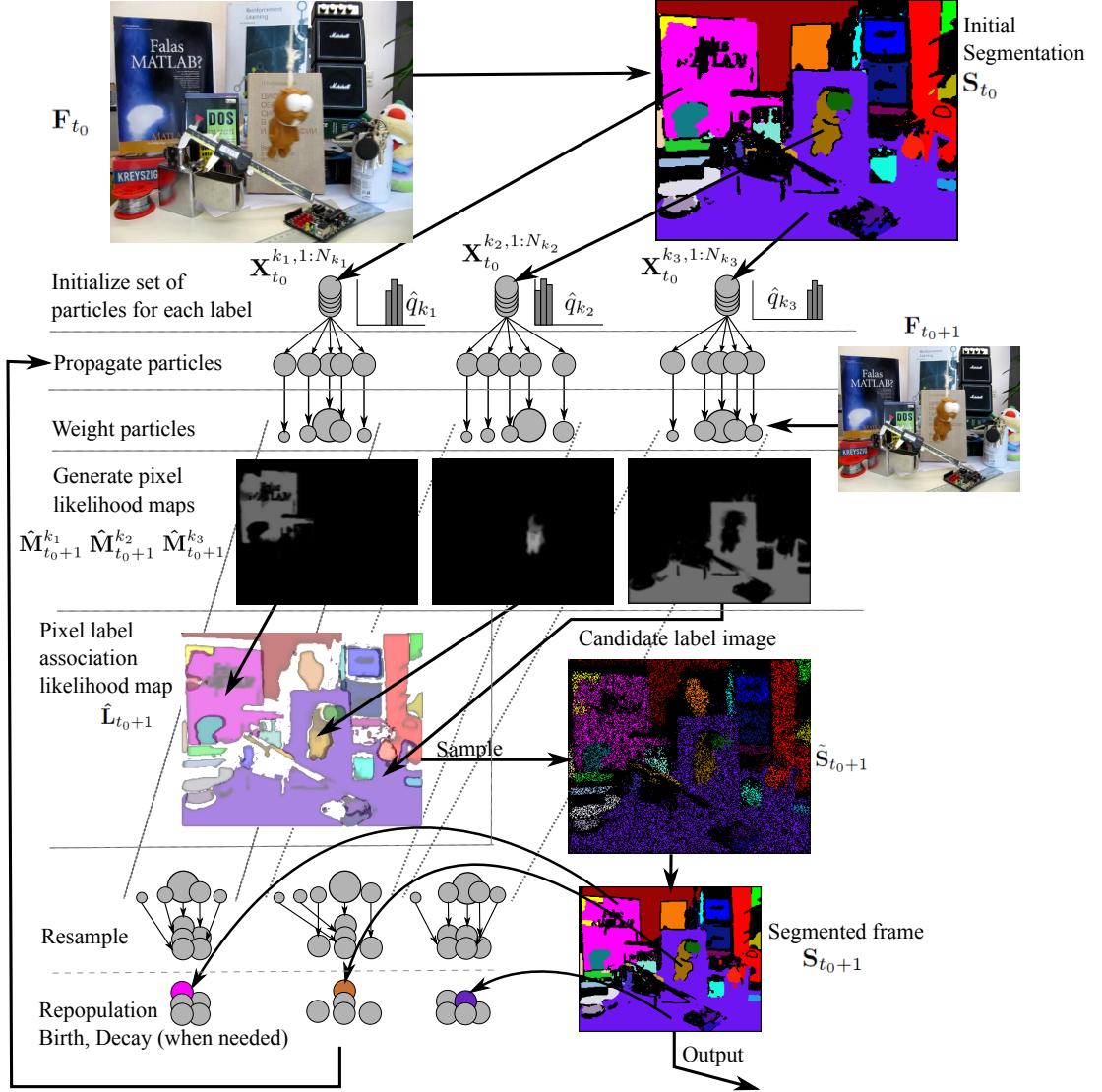
The underlying principle guiding the proposed algorithm is to use predictions from Bayesian filtering to inform segmentation of higher-level temporal object correspondences. It is well known that sequential Bayesian estimation methods perform well in difficult tracking scenarios [30], and, under the Markov assumption, are computationally less demanding than video segmentation techniques such as MHVS [83], which consider many prior frames. Particle filtering is one such method which has been shown to approximate the optimal tracking solution well, even in complex multi-target scenarios with strong nonlinearities [40, 86, 88].

## 2.1 OVERVIEW OF THE ALGORITHM

Before proceeding to discuss elements in detail, we shall first give a brief overview of the algorithm (depicted in Figure 2.1.1). We begin by performing an initial segmentation (using any method) on the first frame  $\mathbf{F}_{t_0}$  to generate an initial set of labels  $\mathbf{S}_{t_0}$ . An initial set of particles is generated for each label, and color histogram features are computed for each particle (as in [69]). Thus each object  $k$  at initial time  $t_0$  is specified by a set of  $N_k$  particles  $\mathbf{X}_{t_0}^{k,1:N_k}$ , each of which contains a representation of the object, specified by a pixel existence map  $\mathbf{M}$ , a reference color histogram  $\mathbf{\hat{h}}$ , a position shift vector  $\mathbf{p}_{t_0}$ , and a velocity vector  $\mathbf{v}_{t_0}$ .

The particles are then propagated in time independently, shifting their existence maps to new regions of the image. These shifted maps are used to generate measured color histograms from the next frame, which are evaluated to determine similarity to the object's reference histogram. The set of particles for each object is then combined to create an overall object pixel likelihood map. The pixel likelihood maps for all objects are then further combined with each other to create a label association likelihood map. In this likelihood map, each pixel is a Probability Distribution Function (PDF) specifying the probability that the original image pixel was generated from an observation of a particular object.

The label association likelihood map is then sampled using a per-pixel selection procedure (as described in Section 2.3.2) to generate a candidate label image,  $\tilde{\mathbf{S}}_{t_0+1}$ . This candidate image is used as the initialization for the Metropolis-Hastings algorithm with annealing of Abramov et al. [2], which updates the labels iteratively until an equilibrium segmented state is reached. The segmentation result,  $\mathbf{S}_{t_0+1}$  is subsequently used to update the set of particles via three mechanisms; birth, decay, and repopulation. Birth is used for new labels in the segmentation output, and consists of initializing a new set of particles. Decay occurs when a label is not found in the segmentation output, and consists of killing a number of the particles of the missing label. The most commonly occurring mechanism, repopulation, occurs for all previously existing object labels which are found. Repopulation rejuvenates the set of particles for



**Figure 2.1.1:** Flow of algorithm for one time step, shown for three labels ( $k_1, k_2$ , and  $k_3$ ). For a description, see Section 2.1.

an object by replacing a number of particles in the set with new particles based on the relaxed segmentation result.

## 2.2 TRACKING OBJECT MASKS

We shall now describe each of the parts of the algorithm given above in further detail, beginning with a description of how we track object masks using particle filters. First we will briefly review the basic principles of sequential Bayesian estimation and particle filtering, and then show how they can be used to predict pixel-level label associations in order to seed a segmentation algorithm.

### 2.2.1 SEQUENTIAL BAYESIAN ESTIMATION

Sequential Bayesian estimation uses a state space representation, in which a state vector  $\mathbf{x}_t$  describes the hidden state of a dynamic system. Bayesian estimation attempts to determine the posterior distribution of the state given all prior observations  $\mathbf{z}$ , i.e.,  $p(\mathbf{x}_t | \mathbf{z}_{1:t})$ . This is accomplished using a two step recursion which first generates a hypothesis of the current state conditioned on the previous state and then performs a Bayes update using the new observation. These steps are known as the prediction and filtering steps, respectively.

The prediction step estimates the current distribution given all prior observations, or

$$p(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}. \quad (2.1)$$

This prediction requires the specification of a stochastic *dynamic model*

$$\mathbf{x}_t = f_t(\mathbf{x}_{t-1}, \mathbf{v}_t), \quad (2.2)$$

where  $\mathbf{v}_t$  is the process noise, which characterizes the state transition density  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ . The dynamic model takes advantage of knowledge of the system to generate reliable predictions of how the state evolves.

The filtering step uses Bayes rule to update the predicted density by conditioning it on the new observation  $\mathbf{z}_t$ :

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}) = \frac{p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{z}_{1:t-1})}{p(\mathbf{z}_t | \mathbf{z}_{1:t-1})}. \quad (2.3)$$

This requires the specification of an observation, or measurement, model

$$\mathbf{z}_t = h_t(\mathbf{x}_t, \mathbf{w}_t), \quad (2.4)$$

where  $\mathbf{w}_t$  is the measurement noise, which characterizes the observation density  $p(\mathbf{z}_t | \mathbf{x}_t)$ . Once the filtered, or posterior distribution is determined, an estimate of the state can be made using a variety of techniques (e.g., maximum a-posteriori, mean-shift).

## DYNAMIC MODEL

In our method, the state of a particle consists of four elements; the pixel existence map  $\mathbf{M}$ , a reference color histogram  $\hat{q}$ , a position shift vector  $\mathbf{p}$ , and a velocity vector  $\mathbf{v}_t$ . Of these, only the position shift and velocity evolve over time, so we adopt the state vector

$$\mathbf{x}_t = [p_x v_x p_y v_y]^T, \quad (2.5)$$

where  $(p_x, p_y)$  denotes the accumulated shift of the pixel existence map in the image plane, and  $(v_x, v_y)$  the map velocity in the image plane. Motion is modeled using a constant velocity

model in discrete time with uniform sampling period  $T$ , giving the dynamic model

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{v}_t, \quad (2.6)$$

where

$$\mathbf{A} = \begin{bmatrix} 1 & T & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & T \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (2.7)$$

and noise  $\mathbf{v}_t$  is assumed to be zero mean Gaussian with fixed covariance.

### MEASUREMENT MODEL

In our method measurements are taken by calculating a color histogram,  $q_t$  for the region lying within the shifted pixel existence map  $\mathbf{M}$ . That is, for particle  $n$  of object  $k$ ,

$$q_t^{k,n} = \text{hist}(\mathbf{F}_t \cap \mathbf{M}_t^{k,n}). \quad (2.8)$$

Color histograms are three dimensional, with 8 bins for each of the color components hue, saturation, and value. As in [69], a Gaussian density is used for the observation density  $p(\mathbf{z}_t | \mathbf{x}_t)$ , that is

$$p(\mathbf{z}_t | \mathbf{x}_t) = \frac{1}{\sqrt{2\pi}\sigma} \exp -\frac{\Delta(\hat{q}, q_t)^2}{2\sigma^2}, \quad (2.9)$$

where  $\Delta(\hat{q}, q_t)$  is the Bhattacharyya distance (as proposed in [29]) between the reference histogram  $\hat{q}$  for the particle and the measured histogram  $q_t$  for time  $t$ . The Bhattacharyya distance is a standard measure of similarity between discrete probability distributions, and is defined as

$$\Delta(\hat{q}, q_t) = \sqrt{1 - \sum \sqrt{\hat{q}q_t}}. \quad (2.10)$$

#### 2.2.2 PARALLEL PARTICLE FILTERS

Except in special cases (e.g., Kalman Filter), closed-form solutions to Equations (2.1) and (2.3) are not available. Particle Filters are a Monte-Carlo method designed to approximate the posterior distribution with a weighted set of random samples. There are many excellent descriptions of the mechanics of particle filtering available (such as [30]), so we shall avoid presenting them here, and proceed directly to presenting the details of our algorithm.

The predictive portion of the method uses multiple Sequential Importance Resampling (SIR) filters in parallel to track multiple targets (labels) simultaneously. At this stage in the algorithm targets are assumed independent and interaction between labels is therefore not considered (interaction is accounted for later, as described in Section 2.3). Particles are first propagated using the constant velocity dynamics model, and their predicted existence maps

$\tilde{\mathbf{M}}^{k,n}$  are used to generate a measured histogram,  $q_t$ . Particles are weighted based on (2.9), and then normalized as a set for each label  $k$ . Systematic resampling is used to prevent particle degeneracy, due to its speed and good empirical performance [30].

The resulting distributions from the weighting procedure are used to generate object pixel likelihood maps for each label,  $\hat{\mathbf{M}}_{t+1}^k$ , which are then combined into the label association likelihood map  $\hat{\mathbf{L}}_t$ , as described in Section 2.3. A realization of this likelihood map can then be relaxed to produce a final segmented output,  $\mathbf{S}_t$ .

### 2.2.3 PARTICLE BIRTH, REPOPULATION, & DECAY.

One key improvement of the proposed algorithm over prior particle filtering methods is its use of the segmentation result  $\mathbf{S}_t$  to update the particle sets. This allows the creation of new targets, adaptation to changing target appearance, and gradual elimination of targets which are no longer observed. This is accomplished via three mechanisms, which we term, respectively, birth, repopulation, and decay.

Birth occurs when a label which has not existed previously is found in the segmentation output  $\mathbf{S}_t$ , or more formally  $\{k \notin \mathbf{S}_{1:t-1}, k \in \mathbf{S}_t\}$ . It consists of generating a set of particles  $\mathbf{X}^k$  for the new label using  $\mathbf{S}_t$  to initialize an existence map  $\mathbf{M}_t^k$  and  $\{\mathbf{F}_t \cap \mathbf{M}_t^k\}$  to calculate a reference color histogram  $\hat{q}_t^k$ .

Repopulation is a key component of the algorithm, as it allows the pixel likelihood map for an object,  $\hat{\mathbf{M}}^k$ , to adapt over time to the changing appearance of the object. Every iteration, all previously existing object labels which are found in  $\mathbf{S}_t$  are repopulated by replacing some particles in the set with particles generated from  $\mathbf{S}_t$  and  $\mathbf{F}_t$ . Particles are chosen for replacement using stratified sampling, at a rate specified by parameter  $\lambda_r$ . The repopulation mechanism gradually modifies the object "model" through the addition of particles which have an updated existence map and color histogram (coming from the segmentation result). We use the term model here loosely, since there is in actuality no explicit model for any of the objects - merely a pixel likelihood map generated at each time step from the objects constituent particles and the current image frame.

Stratified replacement and relatively low repopulation rates are used to help keep the influence of erroneous hypotheses to a minimum, but as with any adaptive method, they can occasionally lead the tracker astray. Replacement of particles, rather than updating of a central model, helps to reduce this problem, since a few erroneous particles will generally not completely derail the algorithm. Nevertheless, future work could investigate strategies that allow pruning of unlikely hypotheses without negatively affecting occlusion handling.

Decay occurs when a label is not found in the segmentation output,  $k \notin \mathbf{S}_t$ . Particles are selected from  $k$  using random sampling, at a rate determined by the decay rate  $\lambda_d$ , and are pruned; they are no longer considered when filtering  $k$ . This reduces the number of active particles for the label in the next iteration,  $N_{t+1}^k$ , by approximately  $\lambda_d N_t^k$ . If the number of active particles for a label falls below a certain threshold,  $N_{min}$ , then the set of particles for

the label is deleted, and the object is no longer tracked. If a label which was being decayed is observed again, i.e.,  $\{k \notin \mathbf{S}_{t-1}, k \in \mathbf{S}_t\}$ , then the label is revived by replacing particles which had been killed with new particles, which are initialized as in the repopulation step.

### 2.3 EXTRACTING A DENSE IMAGE LABELING

The middle portion of Figure 2.1.1 depicts how the candidate label image,  $\tilde{\mathbf{S}}_t$ , is generated. The candidate label image is a summary of the accumulated knowledge of the particle filters; it is a prediction of what the segmented scene should look like. That is to say, it is a pixel-wise realization of the label association likelihood map  $\hat{\mathbf{L}}_t$ , which is constructed by combining the object pixel likelihood maps (which approximate the posteriors of the particle sets).  $\tilde{\mathbf{S}}_t$  is the seed of the segmentation kernel, which uses pixel values from  $\mathbf{F}_t$  to perform the relaxation process and generate a dense label image. In this section we will describe the process of generating the object pixel and label association likelihood maps, and then explain how the predictive loop allows occlusion handling without explicit object relationships or depth modeling.

#### 2.3.1 OBJECT PIXEL LIKELIHOOD MAPS.

The object pixel likelihood map for a particular object  $k$  is the weighted sum of the pixel existence maps of all of its labels,

$$\hat{\mathbf{M}}_t^k = \sum_{n=1}^{N_k} w_t^{k,n} \mathbf{M}^{k,n}. \quad (2.11)$$

Because the weights have been normalized, the pixel values in  $\hat{\mathbf{M}}_t^k$  will be in the range  $[0, 1]$ . High pixel values will occur in regions which are present in the existence maps of highly weighted particles, or alternatively, are present in many particles with average weight.

#### 2.3.2 LABEL ASSOCIATION LIKELIHOOD MAP.

The label association likelihood map  $\hat{\mathbf{L}}_t$  is a combination of all the object pixel likelihood maps, such that each pixel contains a discrete probability distribution giving the likelihood of the pixel belonging to a certain label. Additionally, a likelihood,  $p_0$ , for the pixel belonging to no label is inserted to allow pixels where no label has high likelihood to remain unlabeled in  $\tilde{\mathbf{S}}_t$ . More formally,

$$\hat{\mathbf{L}}_t = \bigcup_{n=1}^K \hat{\mathbf{M}}_t^n + p_0. \quad (2.12)$$

Each pixel of  $\hat{\mathbf{L}}_t$  is then normalized, such that the sum of the discrete probabilities sums to one. The candidate label image can then be generated by taking a realization of  $\hat{\mathbf{L}}_t$  to select pixel label values. Examples of the result of this process,  $\tilde{\mathbf{S}}_t$ , can be seen in Figures 2.1.1 and 2.6.1.

## 2.4 OCCLUSION HANDLING.

Occlusion relationships are handled naturally, since foreground objects will tend to have a strong peak in their weight distribution, corresponding to those particles which align properly with  $\mathbf{F}_t$ . Objects they occlude will have a flat particle weight distribution, since there will exist no shifted existence map which contains a color distribution which matches the reference histogram. This is due to the fact that the occluding objects and objects surrounding the occluded object have color distributions which differ from the occluded object. Let us assume foreground object  $j$  is contained by occluded object  $k$ , that is

$$\mathbf{M}_t^{j,n} \subset \mathbf{M}_t^{k,n}. \quad (2.13)$$

We also assume that the number of particles is sufficiently large such that

$$\exists \mathbf{M}_t^{j,n} \in \mathbf{M}_t^j : hist(\mathbf{F}_t \cap \mathbf{M}_t^{j,n}) \approx \hat{q}^{j,n}. \quad (2.14)$$

If  $hist(\mathbf{F}_t \cap \mathbf{M}_t^{k,n}) \neq hist(\mathbf{F}_t \cap \mathbf{M}_t^{j,n})$ , that is, the objects have different color distributions, then from (2.13) and (2.14), it follows that <sup>1</sup>

$$\nexists \mathbf{M}_t^{k,n} \in \mathbf{M}_t^k : hist(\mathbf{F}_t \cap \mathbf{M}_t^{k,n}) \approx \hat{q}^{k,n} \quad (2.15)$$

and therefore that

$$\begin{aligned} min_{1:N_j} \{ \Delta(\hat{q}^{j,n}, hist(\mathbf{F}_t \cap \mathbf{M}_t^{j,n})) \} < \\ min_{1:N_k} \{ \Delta(\hat{q}^{k,n}, hist(\mathbf{F}_t \cap \mathbf{M}_t^{k,n})) \} \end{aligned} \quad (2.16)$$

and thus

$$max_{1:N_j} \{ w_t^{j,n} \} > max_{1:N_k} \{ w_t^{k,n} \}. \quad (2.17)$$

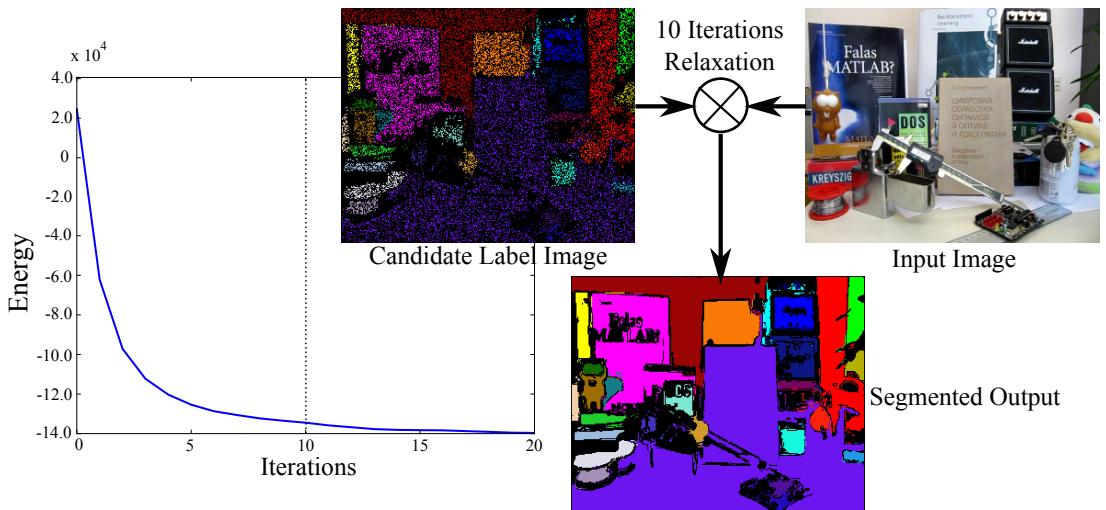
This means that in the label association likelihood map  $\hat{\mathbf{L}}_t$ , the occluding object will have a higher likelihood than the occluded. The candidate label image,  $\tilde{\mathbf{S}}_t$  will therefore tend to favor occluding object labels, which will dominate the occluded object label during the segmentation relaxation process.

## 2.5 SEGMENTATION USING SUPERPARAMAGNETIC CLUSTERING

To adjust the candidate label image  $\tilde{\mathbf{S}}_t$  to the current frame  $\mathbf{F}_t$ , we use a real-time image segmentation algorithm based on superparamagnetic clustering of data [17]. The method of superparamagnetic clustering represents an input image being segmented by a Potts model, with

---

<sup>1</sup>This also assumes that the areas surrounding the occluded object also have different color distributions.



**Figure 2.5.1:** The relaxation process causes the energy of the label image to converge after few iterations (outcome after 10 iterations shown here). This results in efficient calculation of an accurate and temporally coherent segmentation.

pixel color vectors arranged on the sites of a two-dimensional (2D) lattice, where each pixel is featured by an additional variable, called a spin. This allows the segmentation problem to be formulated as a minimization problem which seeks to find the equilibrium states of the energy function in the superparamagnetic phase. In this equilibrium state regions of aligned spins co-exist and correspond to a natural partition of the image data [17]. Since every found segment carries a spin variable which is unique within the whole image, the terms *spin* and *label* are equivalent here. The equilibrium states are found by the use of the highly parallel Metropolis algorithm with a simulated annealing, called *relaxation process*, implemented on a Graphics Processing Unit (GPU) [2]. In this work, the relaxation process adjusts the predicted candidate label image to the current frame.

Superparamagnetic clustering of data was chosen due to its flexibility in allowing the use of any initialization state; there are no particular requirements to the initial states of spin variables. The closer the initial states are to the equilibrium, the less time the Metropolis algorithm needs to converge. This property makes it possible to achieve temporal coherency in the segmentation of temporally adjacent frames by using the sparse label configuration taken from the candidate label image for the spin initialization of the current frame. A final (dense) segmentation result is obtained within a small number of Metropolis updates. Conventional segmentation methods do not generally have this property and cannot turn a sparse segmentation prediction into dense final segments which preserve temporal coherence. Moreover, since the method can directly use sparse predictions as the seed of the segmentation kernel, we can avoid the costly and error-prone block-matching procedure required to find label correspondences in other work, such as in Brendel and Todorovic [21] or Hedau et al. [39]. Figure 2.5.1 illustrates the relaxation process, and the convergence after a small number of iterations.

**Table 2.6.1:** PROST dataset benchmark results. The top table gives average pixel error (lower is better), and the bottom table gives PASCAL based scores (higher is better). Our scores are listed under “HybridPF”. We compare favorably in three of the sequences, and fail on the “box” sequence due to our unsupervised initialization of objects to track.

Sequence	PROST	MIL	Frag	ORF	HybridPF
Lemming	25.1	14.9	82.8	166.3	19.8
Box	13.0	104.6	57.4	145.4	114.1
Liquor	21.5	165.1	30.7	67.3	25.5
Board	39.0	51.2	90.1	154.5	30.9
<hr/>					
Lemming	70.5	83.6	54.9	17.2	73.9
Box	90.6	24.5	61.4	28.3	7.5
Liquor	85.4	20.6	79.9	53.6	54.2
Board	75.0	67.9	67.9	10.0	71.4

## 2.6 EXPERIMENTAL RESULTS

In order to evaluate performance, we compare our method to the state of the art on several challenging video tracking benchmark sequences which are available online<sup>2</sup>. It should be noted that, as opposed to the other tracking algorithms, we do not pre-select a region to track, and track fully deforming object masks (rather than a rectangle). Additionally, we employ no learned or a-priori specified models, use 50 particles per label, and only have two parameters; the repopulation and decay rates  $\lambda$ , and  $\lambda_d$ , which were both held constant at 0.05 throughout testing. Results are compared to the PROST [76], MilTrack [14], FragTrack [6], and ORF [75] tracking algorithms. Further details concerning the parameters used for the above algorithms in the benchmarking can be found in [76].

We shall not evaluate the visual quality of segmentation results here for a couple of reasons. First, detailed evaluation of the visual quality of super-paramagnetic clustering has been presented in [2] in great detail. The visual quality of the segmentation results obtained from this work do not differ significantly from these results, with the exception of labels having continuity through occlusions. Secondly, it is directly acknowledged in other VOS work that the methods fail under partial [55, 64] or full [21, 83] occlusions. As such, comparing performance to other VOS methods is somewhat unreasonable. Rather, the better comparison is to the state of the art in tracking methods, which attempt to handle full and partial occlusions.

In order to compare with the other methods, we needed to output a tracking rectangle for each frame. To do this, once the sequence was segmented, we found the segment which corresponded to the object to track in the first frame, and then took the bounding-box which contained it in each frame as the tracking rectangle. This bounding-box was then compared

---

<sup>2</sup><http://www.GPU4Vision.org>

to ground-truth using two measures; Euclidean distance and the PASCAL-challenge based score proposed in [76]. The latter compares the area of intersection of the ground truth and tracked box with the union of the same. When this is greater than 0.5, the object is considered successfully tracked. Table 2.6.1 gives our results, as well as the results for the other methods.

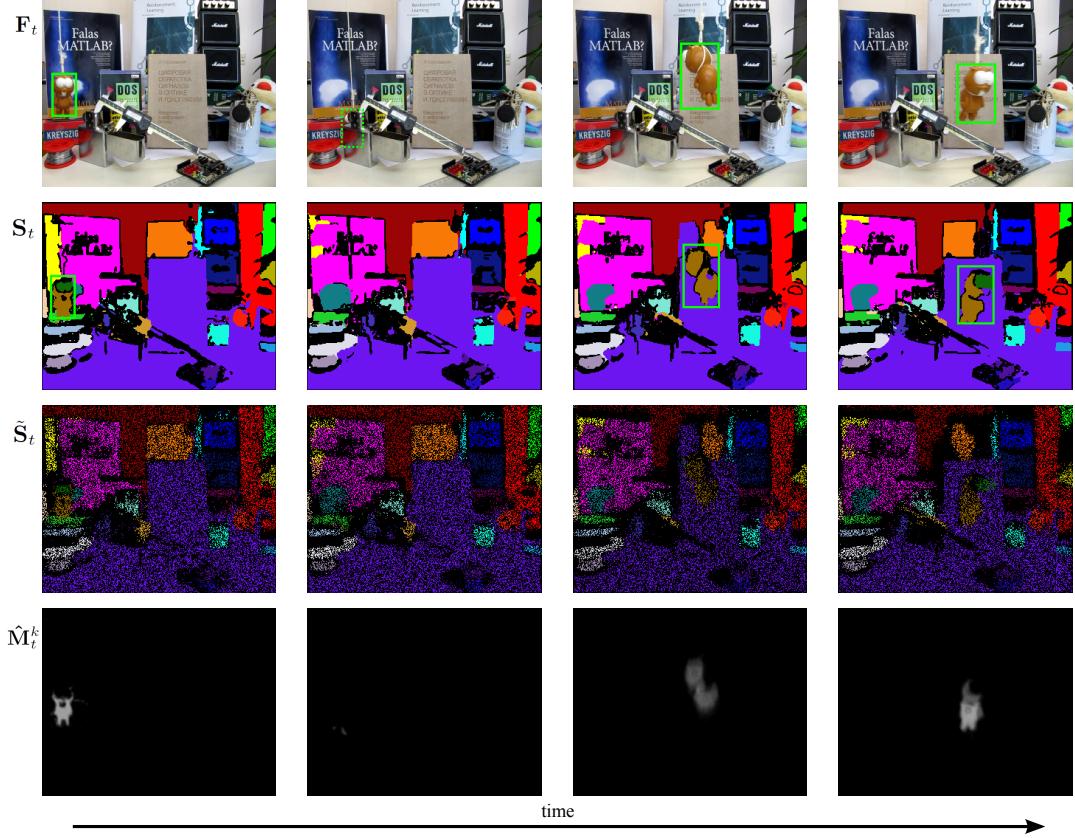
Testing showed that, when certain assumptions hold, our algorithm performs on par with, and in some cases outperforms, state of the art tracking algorithms. This is the case for the *liquor*, *lemming*, and *board* sequences. In the *lemming* sequence, frames of which are shown in Figure 2.6.1, our algorithm outperforms the other methods in cases of occlusion, especially when the tracked object is fully occluded. While other methods offer false positives and erroneous tracks, our method decays the label for the object and avoids proposing incorrect tracking solutions. In the *liquor* sequence, our algorithm adapts to the changing appearance (size, shape, and color) of the tracked bottle, allowing it to maintain performance on par with the other algorithms, in spite of the difficulties of segmenting transparent objects. In the *board* sequence, our method successfully adapts to the rapidly changing appearance of the tracked board as it rotates, allowing it to maintain an accurate track and outperform the other methods.

In addition to showing the strengths of our method, a weakness was also highlighted by the benchmark sequences. The *box* sequence demonstrated the limitations of using unsupervised color-based segmentation to initialize the objects to track. In the sequence, the object to track contains strong color differences, which are segmented into different initial regions. As the object moves around, the particles for these regions are attracted to other objects it passes over which have similar color.

## 2.7 DISCUSSION

In this chapter we presented a new method for performing on-line, dense, unsupervised video segmentation which uses tracking as the basis for segmentation. We have given results which show that the method is able to resolve occlusion relations between objects without explicitly modeling them, and can maintain consistent labels for objects, even when they leave and re-enter the field of view. Additionally, we have shown that the method is able to adapt to rapidly changing appearance of tracked objects, producing consistent segmentations over lengthy video sequences. A GPU version of the algorithm has been developed that can achieve near real-time levels 10 fps at 640x480 resolution on an i7 standard desktop. The algorithm has significant advantages over other VOS methods, in particular when it comes to occlusion handling and speed.

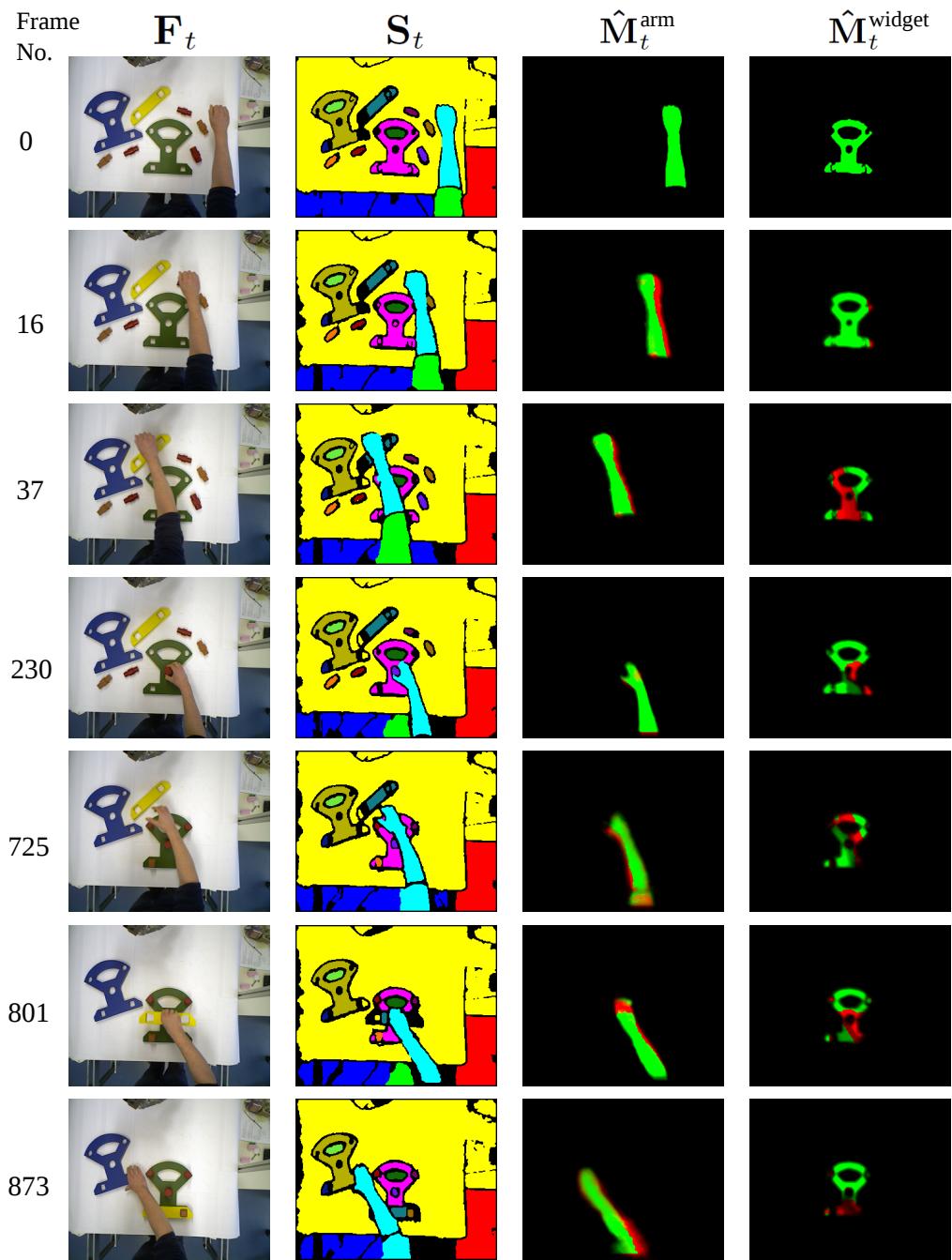
Unfortunately, we found that there were significant instabilities in the algorithm that compounded over time. While these instabilities did not show up in the short sequences used to benchmark, if the algorithm was run for several minutes, it would tend to decay into a state where one or two large segments dominate the others. The other tracking algorithms did not



**Figure 2.6.1:** Output frames from the *lemming* sequence, in which a target is completely occluded (for  $\sim 20$  frames, second column) and changes significantly in appearance. The object which is tracked for comparison to other algorithms is highlighted with a green box.  $\mathbf{F}_t$ -Original frames from movie.  $\mathbf{S}_t$ -The output of the segmentation algorithm.  $\tilde{\mathbf{S}}_t$ -The candidate label image constructed by taking a random draw from  $\hat{\mathbf{L}}_t$ , the label association likelihood map.  $\hat{\mathbf{M}}_t^k$ -The overall object pixel likelihood map for the lemming label, created by combining the set of particles for the label. Intensity represents the sum of the normalized weights of the set of particles.

suffer from this, as they only tracked single targets. We address this problem of instability in subsequent work by reinitializing the segments periodically and finding a solution to the data-association problem.

Another disadvantage of the method is its vulnerability to situations where objects have similar color distributions. This is an inherent flaw of standard (i.e. color-based) video tracking and segmentation systems - they are unable to use 3D shape to resolve ambiguities. As there are many situations where color is not a useful feature for resolving objects, one must reason about real-world geometry in order to reliably track and segment objects. While it is theoretically possible to infer 3D geometry without depth information, the inference tends to be noisy and error-prone, even for the human visual cortex (consider how easy it is to fool one's depth-reasoning with one eye closed). As such, we decided to progress from standard video to RGB-D sensors which provide measured depth information for each pixel, allowing us to incorporate 3D geometry directly into our measurement function and segmentation features. In subsequent Chapters we shall investigate how to take advantage of depth information



**Figure 2.6.2:** Results of segmentation on Cranfield Benchmark Sequence. Green masks show observed pixels, while red masks show occluded pixels which are believed to belong to the object.

while remaining efficient, and how to handle point cloud data which lacks the rigid lattice of image data. With this in mind, in the next Chapter we establish a graph-based framework for efficiently representing streaming unstructured point cloud data.



*The world is a construct of our sensations, perceptions, memories. It is convenient to regard it as existing objectively on its own. But it certainly does not become manifest by its mere existence.*

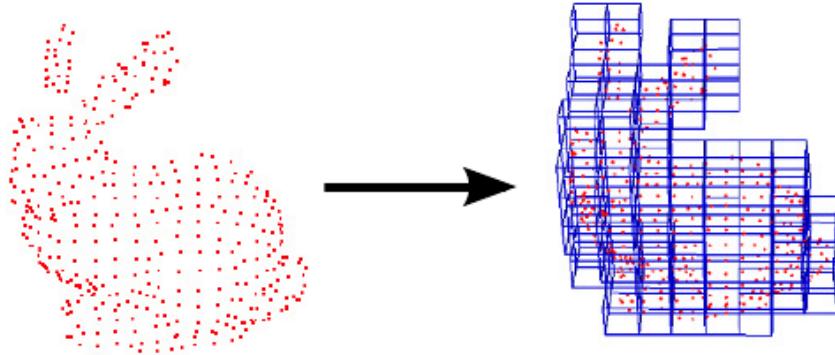
Erwin Schödinger

# 3

## Patch-based Perceptual World Model

THE WIDESPREAD AVAILABILITY of cheap 3D sensors has had a profound impact on the world of computer vision. Where before researchers needed to find heuristic tricks or complex algorithms with which to infer an artificial three dimensional interpretation from a two dimensional image, the new sensors allow direct observation (albeit, noisily) of 3D data. This has allowed direct progression to high-level concepts and rules which the human mind uses when first learning to understand the real world. This is a completely different approach than trying to mimic the behavior of the mind when it is adapting those rules (learned from a life-time of 3D stereo data) in order to interpret some new 2D image. In other words, working within the 3D representation directly allows us to side-step the problem of needing to imitate the complex machinery[68] the mind uses to construct an internal representation of the world.

In this chapter we shall present our work in creating such a full 3D artificial world model which can be used for efficient higher level semantic understanding of both single frames and video. While we do not claim that the model proposed in this Chapter bears direct similarity to the one used internally in the visual cortex, we have found its use generally advantageous over the 2D projective representation. Indeed, we suggest that the concept of “empty space” which is encoded implicitly in our sparse voxel model is an extremely useful and important notion. Moreover, the model is able to succinctly and unambiguously express spatial relationships as a 2D model cannot.



**Figure 3.1.1:** Illustration of Voxelization. On the left we have a point cloud of the “Stanford Bunny”. This cloud is inserted into the voxel grid shown on the right, where all points falling within one grid unit, or voxel, are combined. From <http://www.pointclouds.org/>

### 3.1 PRE-PROCESSING OF POINT CLOUD DATA

Our model begins with point clouds, relying on the general framework set up in the Point Cloud Library <sup>1</sup>, which we have both made use-of and contributed-to as part of this work. Point clouds are a useful way of representing the data obtained from RGB-D sensors, where pixel coordinates and depth value from the RGB-D pair are transformed into an  $(x, y, z)$  point in continuous real-world space, with the RGB information for the pixel attached to this point. Before continuing, we shall briefly introduce two important pre-processing steps which are used throughout the rest of this work. The first down-samples the continuous point cloud space onto a discrete grid, while the second pre-computes an adjacency graph for this grid.

#### 3.1.1 VOXELIZATION

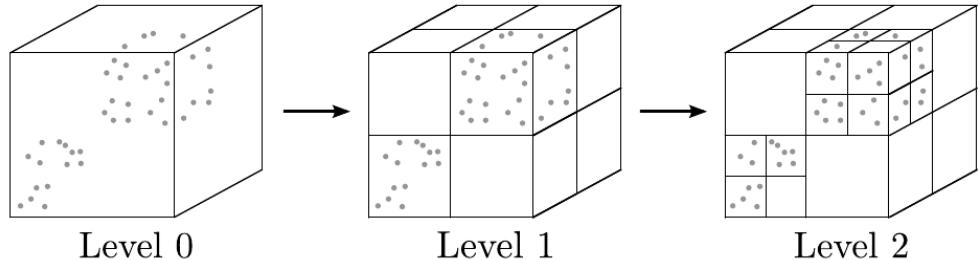
The resolution of a standard RGB-D camera such as the Kinect is 640x480 pixels, yielding about 300,000 points per frame. While for static image segmentation this might be an acceptable amount of data, for video segmentation it is simply too much data to process directly in reasonable run times (on standard hardware). Because of this, a common pre-processing step is to down-sample point clouds using a *voxel-grid* filter, a process known as *voxelization*.

#### 3.1.2 OCTREE ADJACENCY GRAPH

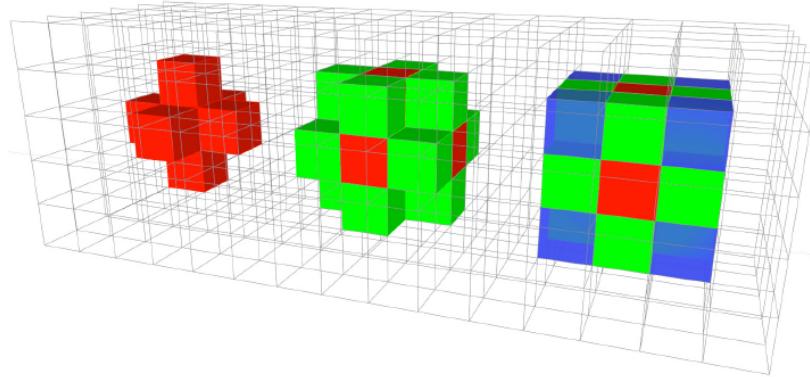
In order to increase computational efficiency, we have developed an adjacency octree which maintains neighbor information within the octree leaves (i.e., the voxels). Adjacency is a key element of many methods, especially region growing algorithms, as it ensures that labels do

---

<sup>1</sup><http://www.pointclouds.org/>



**Figure 3.1.2:** Use of an octree for voxelization. The points are grouped into voxels by recursively subdividing the bounding box into its eight constituent octants. This recursion terminates when the box size has edge length equal to the voxel leaf size  $R_{voxel}$ .



**Figure 3.1.3:** Adjacency in a 3d voxel grid. The 6-, 18-, and 26-neighborhoods share a face, edge, and vertex, respectively.

not cross object boundaries which are disconnected in space. There are three definitions of adjacency in a voxelized 3D space; 6-, 18-, or 26-adjacent. These share a face, faces or edges, and faces, edges, or vertices, respectively. In this work we use 26-adjacency exclusively, as the other lesser adjacencies might miss connections when surfaces are placed in certain configurations relative to the camera plane.

Throughout the rest of this work, we shall deal exclusively with voxels, rather than points, and shall always use our adjacency octree. As voxelization is a necessary pre-processing step for all of the algorithms we shall subsequently discuss, it can be assumed that adjacency information is always available in constant time. This is especially important for the clustering algorithm we introduce next.

## 3.2 GEOMETRICALLY CONSTRAINED SUPERVOXELS

In this Section we present Voxel Cloud Connectivity Segmentation (VCCS), a new method for generating superpixels and supervoxels from 3D point cloud data. The supervoxels produced by VCCS adhere to object boundaries better than state-of-the-art methods while remaining efficient enough to use in online applications. VCCS uses a variant of k-means clus-

tering for generating its labeling of points, with two important constraints:

1. The seeding of supervoxel clusters is done by partitioning 3D space, rather than the projected image plane. This ensures that supervoxels are evenly distributed according to the geometry of the scene.
2. The iterative clustering algorithm enforces strict spatial connectivity of occupied voxels when considering points for clusters. This means that supervoxels strictly cannot flow across boundaries which are disjoint in 3D space, even though they are connected in the projected plane.

First, in 3.2.1 we shall describe how supervoxel seeds are generated and filtered, in 3.2.2 the features and distance measure used for clustering, and finally in 3.2.3 how the iterative clustering algorithm enforces spatial connectivity. Unless otherwise noted, all processing is being performed in the 3D voxelized space constructed from one or more RGB+D cameras (or any other source of point-cloud data). Furthermore, because we work exclusively in a voxel-cloud space (rather than the continuous point-cloud space), we shall adopt the following notation to refer to voxel at index  $i$  within voxel-cloud  $V$  of voxel resolution  $r$ :

$$V_r(i) = \mathbf{F}_{1..n}, \quad (3.1)$$

where  $\mathbf{F}$  specifies a feature vector which contains  $n$  point features (e.g. color, location, normals).

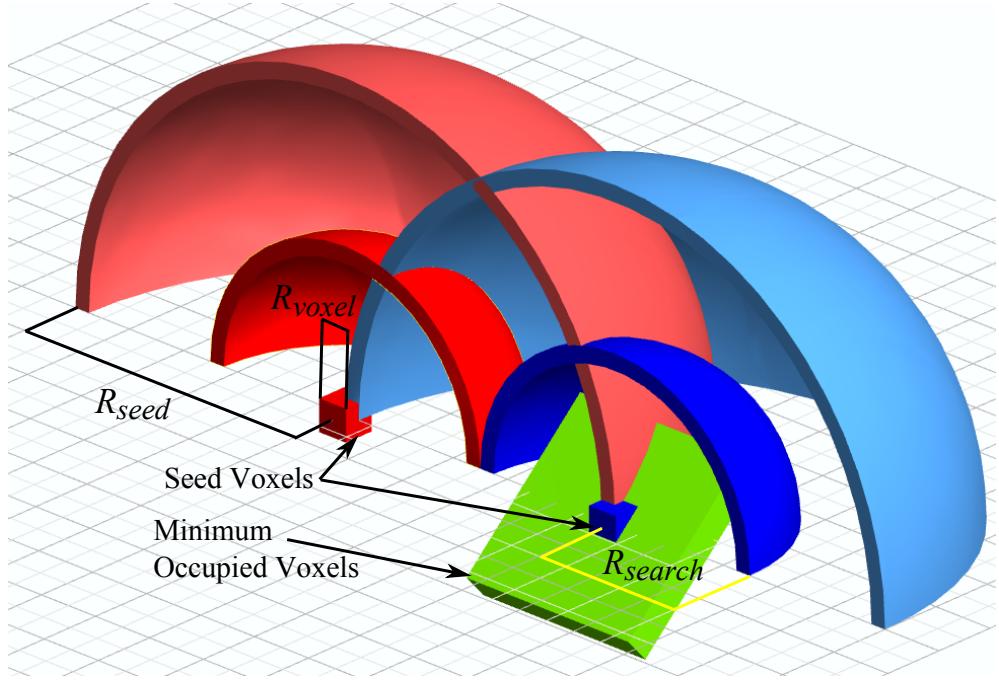
### 3.2.1 SPATIAL CLUSTER SEEDING

The algorithm begins by selecting a number of seed points which will be used to initialize the supervoxels. In order to do this, we first divide the space into a voxelized grid with a chosen resolution  $R_{seed}$ , which is significantly higher than  $R_{voxel}$ . The effect of increasing the seed resolution  $R_{seed}$  can be seen in Figure 3.2.2. Initial candidates for seeding are chosen by selecting the voxel in the cloud nearest to the center of each occupied seeding voxel.

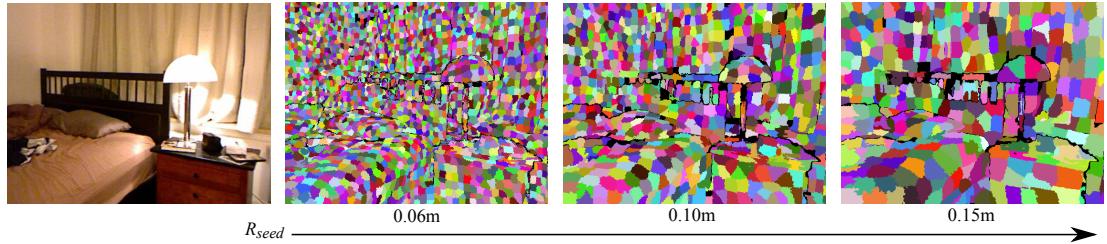
Once we have candidates for seeding, we must filter out seeds caused by noise in the depth image. This means that we must remove seeds which are points isolated in space (which are likely due to noise), while leaving those which exist on surfaces. To do this, we establish a small search radius  $R_{search}$  around each seed, and delete seeds which do not have at least as many voxels as would be occupied by a planar surface intersecting with half of the search volume (this is shown by the green plane in Figure 3.2.1). Once filtered, we shift the remaining seeds to the connected voxel within the search volume which has the smallest gradient in the search volume. Gradient is computed as

$$G(i) = \sum_{k \in V_{adj}} \frac{\| V(i) - V(k) \|_{CIELab}}{N_{adj}}; \quad (3.2)$$

we use sum of distances in CIELAB space from neighboring voxels, requiring us to normal-



**Figure 3.2.1:** Seeding parameters and filtering criteria.  $R_{seed}$  determines the distance between supervoxels, while  $R_{voxel}$  determines the resolution to which the cloud is quantized.  $R_{search}$  is used to determine if there are a sufficient number of occupied voxels to necessitate a seed.



**Figure 3.2.2:** Image segmented using VCCS with seed resolutions of 0.1, 0.15 and 0.2 meters.

ize the gradient measure by number of connected adjacent voxels  $N_{adj}$ . Figure 3.2.1 gives an overview of the different distances and parameters involved in seeding.

Once the seed voxels have been selected, we initialize the supervoxel feature vector by finding the center (in feature space) of the seed voxel and connected neighbors within 2 voxels.

### 3.2.2 CLUSTER FEATURES AND DISTANCE

VCCS supervoxels are clusters in a 39 dimensional space, given as

$$\mathbf{F} = [x, y, z, L, a, b, \text{FPFH}_{1..33}], \quad (3.3)$$

where  $x, y, z$  are spatial coordinates,  $L, a, b$  are color in CIELab space, and  $\text{FPFH}_{1..33}$  are the 33 elements of Fast Point Feature Histograms (FPFH), a local geometrical feature proposed by Rusu et al. [74]. FPFH are pose-invariant features which describe the local surface model properties of points using combinations of their  $k$  nearest neighbors. They are an extension of the older Point Feature Histograms optimized for speed, and have a computational complexity of  $O(n \cdot k)$ .

In order to calculate distances in this space, we must first normalize the spatial component, as distances, and thus their relative importance, will vary depending on the seed resolution  $R_{seed}$ . Similar to the work of Achanta et al., [5] we have limited the search space for each cluster so that it ends at the neighboring cluster centers. This means that we can normalize our spatial distance  $D_s$  using the maximally distant point considered for clustering, which will lie at a distance of  $\sqrt{3}R_{seed}$ . Color distance  $D_c$ , is the euclidean distance in CIELab space, normalized by a constant  $m$ . Distance in FPFH space,  $D_f$ , is calculated using the Histogram Intersection Kernel [15]. This leads us to a equation for normalized distance  $D$ :

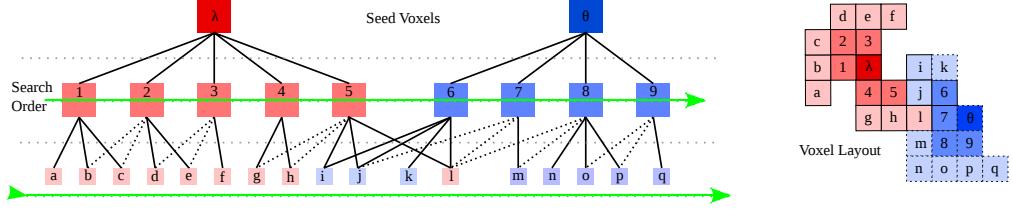
$$D = \sqrt{\frac{\lambda D_c^2}{m^2} + \frac{\mu D_s^2}{3R_{seed}^2} + \varepsilon D_{HiK}^2}, \quad (3.4)$$

where  $\lambda$ ,  $\mu$ , and  $\varepsilon$  control the influence of color, spatial distance, and geometric similarity, respectively, in the clustering. In practice we keep the spatial distance constant relative to the other two so that supervoxels occupy a relatively spherical space, but this is not strictly necessary. For the experiments in this paper we have color weighted equally with geometric similarity.

### 3.2.3 FLOW CONSTRAINED REGION GROWING

Assigning voxels to supervoxels is done iteratively, using a local k-means clustering related to [5, 89], with the significant difference that we consider connectivity and flow when assigning pixels to a cluster. The general process is as follows: beginning at the voxel nearest the cluster center, we flow outward to adjacent voxels and compute the distance from each of these to the supervoxel center using Equation 3.4. If the distance is the smallest this voxel has seen, its label is set, and using the adjacency graph, we add its neighbors which are further from the center to our search queue for this label. We then proceed to the next supervoxel, so that each level outwards from the center is considered at the same time for all supervoxels. We proceed iteratively outwards until we have reached the edge of the search volume for each supervoxel (or have no more neighbors to check).

This amounts to a breadth-first search of the adjacency graph, where we check the same level for all supervoxels before we proceed down the graphs in depth. Importantly, we avoid edges to adjacent voxels which we have already checked this iteration. The search concludes for a supervoxel when we have reached all the leaf nodes of its adjacency graph or none of the



**Figure 3.2.3:** Search order for the flow constrained clustering algorithm (shown in 2D for clarity). Dotted edges in the adjacency graph are not searched, as the nodes have already been added to the search queue.

nodes searched in the current level were set to its label. This search procedure, illustrated in Figure 3.2.3, has two important advantages over existing methods:

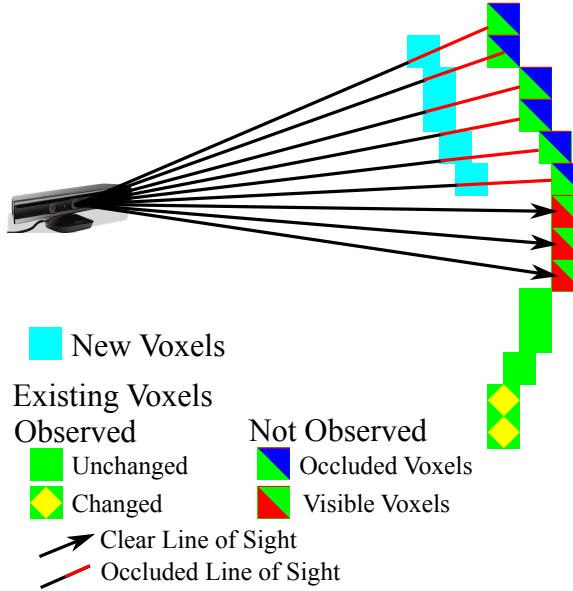
1. Supervoxel labels cannot cross over object boundaries that are not actually touching in 3D space, since we only consider adjacent voxels, and
2. Supervoxel labels will tend to be continuous in 3D space, since labels flow outward from the center of each supervoxel, expanding in space at the same rate.

Once the search of all supervoxel adjacency graphs has concluded, we update the centers of each supervoxel cluster by taking the mean of all its constituents. This is done iteratively; either until the cluster centers stabilize, or for a fixed number of iterations. For this work we found that the supervoxels were stable within a few iterations, and so have simply used five iterations for all presented results.

### 3.3 SEQUENTIAL UPDATE OF PERCEPTUAL MODEL

As an additional consideration, we have developed a scheme for adding new point clouds sequentially (as from a video stream) into an existing supervoxel octree. This is accomplished through a process which classifies voxels in the tree based on their behavior. As a first step, we insert points from the new point cloud into the octree, and initialize new leaves for voxels which did not exist previously. This results in an octree where leaves fall into three possible categories (illustrated in Figure 3.3.1; they are either new, observed, or unobserved in the most recent observation. Handling of new leaves is straightforward; we simply calculate adjacency relations to existing leaves and flag them as unlabeled.

To determine whether a leaf which existed previously has changed, we test the distance between the centroid of the points falling within its voxel (from the new frame) and its previous centroid. This is done in the same feature space used for growing the supervoxels, that is, we test whether the normal, color, and spatial location have varied more than a threshold value. This threshold is set to a relatively low constant value so that it favors false-positives (finding change when there was none), as they do not impact the tracking performance of the algorithm, but only have a slight effect on its run-time. If a leaf is found to have changed, we

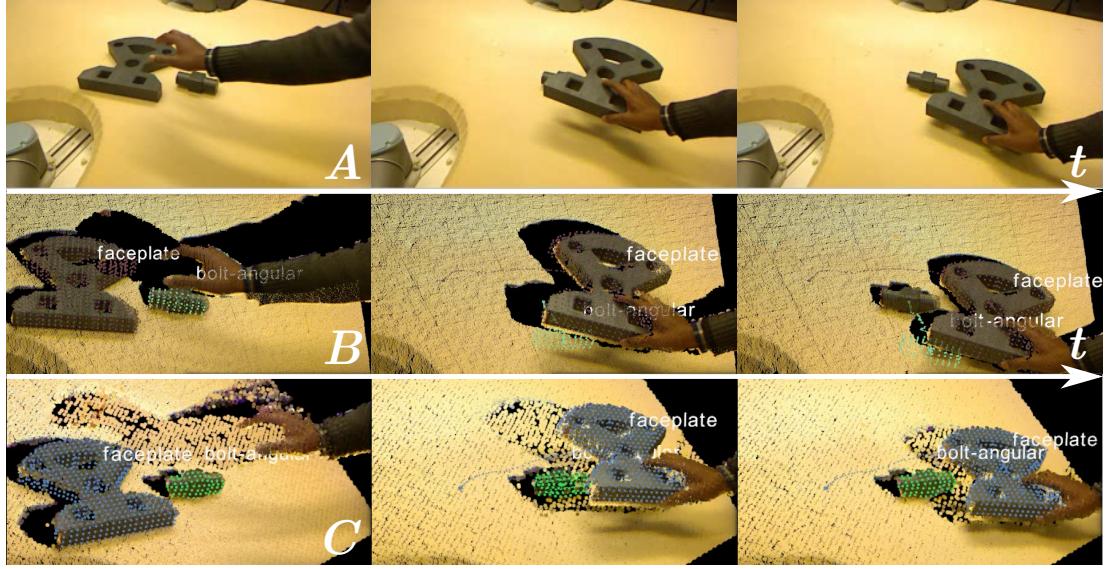


**Figure 3.3.1:** Categorization of voxels based on new frame of data. Voxels fall into three categories, they are either new, observed or not observed in the frame. Furthermore, observed voxels can either have changed or remained the same, while voxels not observed in the frame are either occluded or no longer exist (in which case they should be deleted).

remove its previous labeling. We also perform a global check to see if more than half of a supervoxels support has changed; if so, we completely remove the supervoxels label from all of its constituent voxels.

Finally, we must consider how to handle leaves which were not observed in the inserted point cloud. Rather than simply prune them, we first check if it was possible to observe them from the viewpoint of the sensor which generated the input cloud. This occlusion check can be accomplished efficiently using the octree by determining if any voxels exist between unobserved leaves and the sensor viewpoint. If a clear line of sight exists from the leaf to the camera, it can safely be deleted. Conversely, if the path is obstructed, we "freeze" the leaf, meaning that it will remain constant until it is either observed or passes the line of sight test in a future frame (in which case, it can be safely deleted). This occlusion testing means that tracking of occluded objects is trivial, as occluded voxels remain in the observations which are used for tracking. This procedure results in what we term "voxel-permanence", as it results in voxels persisting through occlusions as seen in Figure 3.3.2.

Once the octree voxels have been updated, we then proceed to update the supervoxels as before. That is, first we generate new seeds in regions of large unlabeled voxels, and then conduct the iterative region growing. This results in new supervoxels in regions which are new or changing, while leaving supervoxels in static and occluded regions unchanged. This reduces the tracking and segmentation problem to finding the best joint association of these new supervoxels with those from the prior time-step.



**Figure 3.3.2:** Example of successful tracking of an object through complete occlusion using the sequentially updated world model and the concept of voxel permanence. Row A shows the original image frames - the bolt becomes occluded by the faceplate and cannot be seen by the sensor. Row B shows tracking failure using the raw 3D data - black “holes” behind the arm and faceplate are due to occlusion. Row C shows our model and tracked output - “holes” are now mostly filled in allowing tracking to succeed.

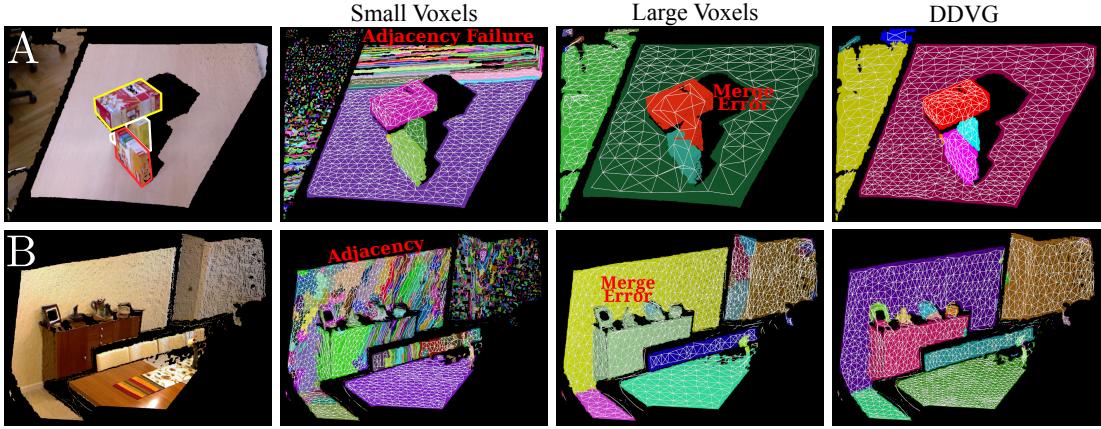
### 3.4 DEPTH DEPENDENT VOXEL GRID

So far we have described the main algorithm for generating supervoxels. Next we will introduce a depth transform which improves supervoxels by addressing the shortcomings of the adjacency octree upon which VCCS depends. As with any system which uses projective geometry, observations from a single RGB-D camera have a significant drawback - the level of detail decreases with increasing distance from the camera. In our case, this manifests as decreasing point density. In addition, the levels of both quantization and noise grow quadratically with distance [47, 80]. The combined effect of quantization and change in point density with depth results in inevitable failure of adjacency computation. At a certain distance (dependent on the voxel size  $R_{voxel}$ ), the sparsity of observed points results in “holes” in the octree, and a break-down of adjacency. This has obvious negative consequences for flow-constrained algorithms such as VCCS which rely on spatial connectivity for clustering.

We compensate for the loss of point density and quantization with increasing depth  $z$  by transforming the points into a skewed space using the transformation  $T : (x, y, z) \rightarrow (x', y', z')$  with

$$x' = x/z, \quad y' = y/z, \quad z' = \log(z) \quad (3.5)$$

The division of the  $x$  and  $y$  coordinates by  $z$  reverses the perspective transformation, equalizing the point density in the  $x$ - $y$ -plane. Transforming the  $z$  coordinate helps to deal with the effects



**Figure 3.4.1:** Two example point clouds (**A,B, left**) showing the need for the DDVG. For better visibility outlines have been drawn around the boxes in A. Using *Small Voxels* objects close to the camera can be segmented, but adjacency breaks down as the depth increases and the point density decreases. Using *Large Voxels* corrects the adjacency graph in the background, but leads to objects being merged in the foreground due to the coarse resolution. Using DDVG, the scale of the voxels gradually increases with distance from the camera – adapting to the increased noise level and lower point density – consequently adjacency is maintained and the segmentation of scenes with large depth variance is possible using fixed parameters.

of depth quantization by compressing points as depth increases. It is easy to show that the transformation has the following property:

$$\frac{\partial x'}{\partial x} = \frac{\partial y'}{\partial y} = \frac{\partial z'}{\partial z} = \frac{1}{z} \quad (3.6)$$

Because the derivatives are equal, the local coordinate frame is stretched equally along all axes by the transformations. The important thing about this property is, that small cubic voxels are still cubic after the transformation. This leaves the geometry of space basically untouched in the foreground (if the voxel size is chosen sufficiently small), while distant voxels are strongly transformed to fill the “empty” space, compensating for reduced point density.

Rather than transforming the clouds back and forth, we instead transform the bins of the octree itself, creating an octree where bin volume (and thus, voxel size) effectively increases with distance from the camera viewpoint. Doing this directly within the octree allows us to determine adjacency as before (neighboring bins), even though distance between neighboring voxels increases with distance from the camera. Figure 3.4.1 illustrates the advantageous effect of this transformation on segmentation.

### 3.5 LOCALLY CONVEX CONNECTED PATCHES

As an example of an application of supervoxels and the adjacency octree, we shall briefly present a segmentation method which breaks a supervoxel adjacency graph into meaningful segments by classifying whether an edge  $e = (\vec{p}_i, \vec{p}_j)$  between two supervoxels is convex or concave. This

classification is based on an *Extended Convexity Criterion (ECC)*, which considers adjacent supervoxels with centroids at the positions  $\vec{x}_1, \vec{x}_2$  and normals  $\vec{n}_1, \vec{n}_2$ . Whether the connection between these is convex or concave can be inferred from the relation of the surface normals to the vector joining their centroids - an overview of the this algorithm is given in Figure 3.5.1.

The angle of the normals to the vector  $\vec{d} = \vec{x}_1 - \vec{x}_2$  joining the centroids can be calculated using the identity for the dot product  $\vec{a} \cdot \vec{b} = |\vec{a}| \cdot |\vec{b}| \cdot \cos(\alpha)$  with  $\alpha = \angle(\vec{a}, \vec{b})$ . For *convex* connections,  $\alpha_1$  is smaller than  $\alpha_2$ . This can be expressed as:

$$\alpha_1 < \alpha_2 \Rightarrow \cos(\alpha_1) - \cos(\alpha_2) > 0 \Leftrightarrow \vec{n}_1 \cdot \hat{d} - \vec{n}_2 \cdot \hat{d} > 0,$$

where  $\hat{d} = \frac{\vec{x}_1 - \vec{x}_2}{\|\vec{x}_1 - \vec{x}_2\|}$ . Similarly, for a *concave* connection we get:

$$\alpha_1 > \alpha_2 \Leftrightarrow \vec{n}_1 \cdot \hat{d} - \vec{n}_2 \cdot \hat{d} < 0.$$

Note that these operations are commutative, thus the choice of which patch is  $\vec{x}_1$ , does not change the result. Also the criterion is still valid if the  $\vec{x}_i$  are displaced, as long as they stay within the surface.

To compensate for noise in the RGB-D data, a bias is introduced to treat concave connections with very similar normals, that is

$$\beta = \angle(\vec{n}_1, \vec{n}_2) = |\alpha_1 - \alpha_2| = \cos^{-1}(\vec{n}_1 \cdot \vec{n}_2) < \beta_{\text{Thresh}},$$

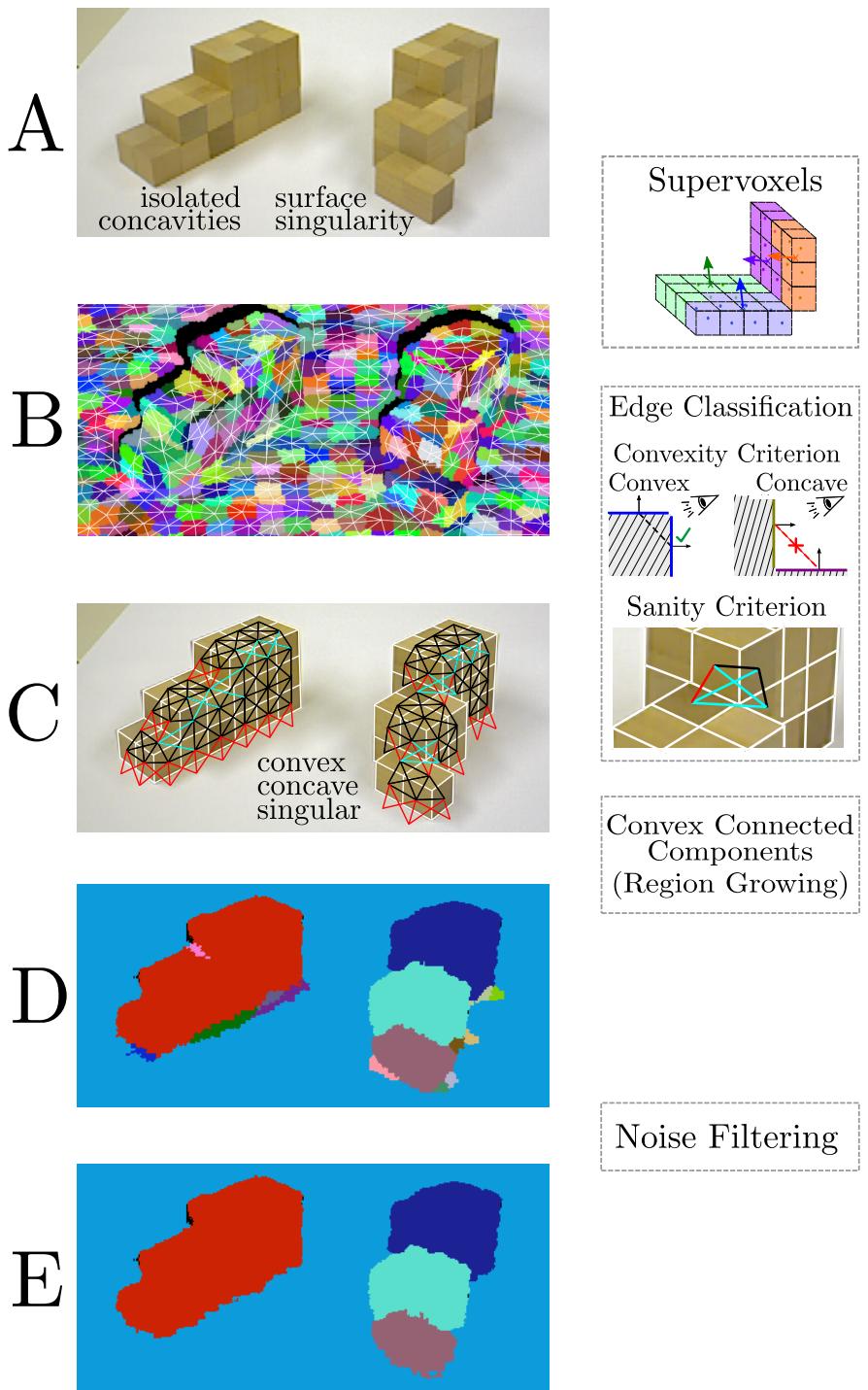
as convex, since those usually represent flat surfaces. Depending on the value of the *concavity tolerance threshold*  $\beta_{\text{Thresh}}$ , concave surfaces with low curvature are seen as convex and thus merged in the segmentation. This behavior may be desired to ignore small concavities. We set:

$$\text{CC}_b(\vec{p}_i, \vec{p}_j) := \begin{cases} \text{true} & (\vec{n}_1 - \vec{n}_2) \cdot \hat{d} > 0 \vee (\beta < \beta_{\text{Thresh}}) \\ \text{false} & \text{otherwise.} \end{cases} \quad (3.7)$$

where the variable  $\text{CC}_b$  defines the *basic convexity criterion*. However, local errors in the feature estimation caused by noise in the data can propagate very easily, potentially leading to errors in the resulting segmentation. This also makes the recognition of small concavities harder, as subtle features are more sensitive to noise. To improve on this we also include neighborhood information in the classification of edges: For a convex edge  $e = (\vec{p}_i, \vec{p}_j)$ , we require that there exists a common neighbor  $\vec{p}_c$  of  $\vec{p}_i$  and  $\vec{p}_j$  that has a convex connection to both.

Thus we define *extended convexity*  $\text{CC}_e$ :

$$\begin{aligned} \text{CC}_e(\vec{p}_i, \vec{p}_j) = & \text{CC}_b(\vec{p}_i, \vec{p}_j) \wedge \text{CC}_b(\vec{p}_i, \vec{p}_c) \\ & \wedge \text{CC}_b(\vec{p}_j, \vec{p}_c) \end{aligned} \quad (3.8)$$



**Figure 3.5.1:** Flow diagram of the segmentation algorithm. **A)** RGB images corresponding to the point clouds of the scene. The red lines show two isolated concavities. The blue box shows an area with a surface singularity. **B)** Supervoxel adjacency graph. **C)** Model depicting the classified graph. Black lines denote convex connections, red lines concave ones and turquoise lines singular connections (those, where two patches are connected only in a single point). **D)** Segmentation result; object labels are shown by different colors. **E)** Final result after noise filtering. The right column illustrates the supervoxel patches and the convexity and sanity criteria used for edge classification.

With extended convexity, more evidence is necessary for a connection to be labeled as convex.

As in VCCS, clusters are found in LCCP using a region growing process: First, an arbitrary seed supervoxel is chosen and labeled. This label is then propagated over the graph with a depth search that is only allowed to grow over convex edges. Once no new supervoxel can be assigned to the segment, we choose a new seed supervoxel that has not been labeled and propagate the new label as before, repeating the process until all supervoxels have been labeled. Note that all of the criteria in LCCP are commutative, so the output of the region growing does not depend on the choice of the seeds.

## 3.6 EXPERIMENTAL RESULTS

### 3.6.1 DATASETS

In the following sections we present quantitative results for VCCS and LCCP. We compare both to state-of-the-art methods on the *NYU Indoor Dataset*[79] and *Object Segmentation Database*[72]. Before giving results, we shall first describe the datasets as well as the procedure for scoring results using 2D ground-truth.

#### OBJECT SEGMENTATION DATABASE (OSD)

The *Object Segmentation Database* (OSD-v0.2) was proposed by Richtsfeld *et al.*[72] in 2012. It consists of 111 cluttered scenes of objects on a table, taken with close proximity to the pictured objects. The scenes contain multiple objects, which have mostly box-like or cylindrical shape, with partial and full occlusions and heavy clutter in 2D as well as 3D. Importantly, most objects in the data set are *simple*, that is, consist of only a single part. This makes the ground-truth data relatively non-ambiguous.

#### NYU INDOOR DATASET (NYU)

The *NYU Indoor Dataset*<sup>2</sup> (NYUv2) from Silberman *et al.*[79] is a large and complex dataset, consisting of 1449 cluttered indoor scenes. The data consists of pairs of aligned RGB and depth images, along with human annotated densely labeled ground truth. The images were captured in diverse indoor scenes, and present many difficulties for segmentation algorithms such as varied illumination and many small similarly colored objects. Examples of typical scenes are shown in Figure 3.6.3. One main difficulty presented by the dataset is that the distance to objects from the camera is quite large in the dataset. This results in significant depth quantization artifacts as well as few data points for many objects. Additionally, depth is often

---

<sup>2</sup>[http://cs.nyu.edu/~silberman/datasets/nyu\\_depth\\_v2.html](http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html)



**Figure 3.6.1:** Example results for scenes from the NYU dataset using unsmoothed depth. Black areas indicate missing depth. Top row: rgb images. Mid. row: segmentation result. Bottom row: ground truth. Parameters A-C:  $R_{voxel} = 0.0075$ ,  $R_{seed} = 0.03$  and  $\beta_{Thresh} = 8^\circ$ . Parameters D-E:  $R_{voxel} = 0.01$ ,  $R_{seed} = 0.04$  and  $\beta_{Thresh} = 10^\circ$  (identical to quantitative results, see Tab. 3.6.2).

missing for extensive portions of many of the images, due to limitations of the Kinect sensor (e.g. reflective, transparent surfaces - windows are especially problematic). Silberman *et al.* attempt to correct for these errors using a hole filling algorithm (*smoothdepth*), which estimates depth for missing areas based on the scheme from Levin *et al.*[53].

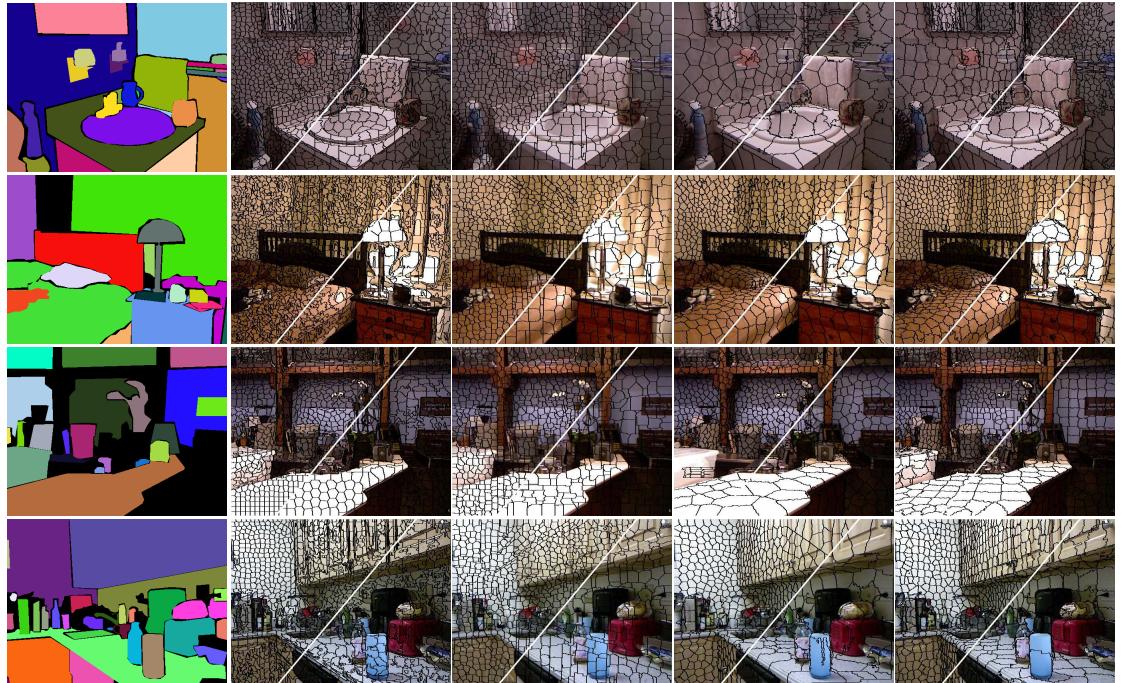
#### RETURNING TO THE PROJECTED PLANE

RGB+D sensors produce what is known as an organized point cloud- a cloud where every point corresponds to a pixel in the original RGB and depth images. When such a cloud is voxelized, it necessarily loses this correspondence, and becomes an unstructured cloud which no longer has any direct relationship back to the 2D projected plane. As such, in order to compare results with existing 2D methods we were forced to devise a scheme to apply supervoxel labels to the original image.

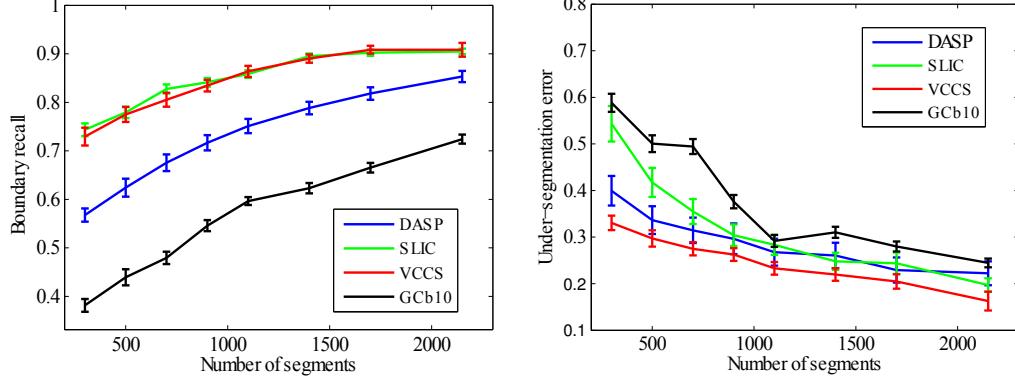
To do this, we take every point in the original organized cloud and search for the nearest voxel in the voxelized representation. Unfortunately, since there are blank areas in the original depth image due to such factors as reflective surfaces, noise, and limited sensor range, this leaves us with some blank areas in the output labeled images. To overcome this, we fill in any large unlabeled areas using the SLIC algorithm. This is not a significant drawback, as the purpose of the algorithm is to form supervoxels in 3D space, not superpixels in the projected plane, and this hole-filling is only needed for comparison purposes. Additionally, the hole filling actually makes our results worse, since it does not consider depth, and therefore tends to bleed over some object boundaries that were correctly maintained in the supervoxel representation. An example of what the resulting segments look like before and after this procedure



**Figure 3.6.2:** Example of hole-filling for images after returning from voxel-cloud to the projected image plane. Depth data, shown in the top left, has holes in it, shown as dark blue areas (here, due to the lamp interfering with the Kinect). The resulting supervoxels do not cover these holes as shown in the bottom left, since the cloud has no points in them. To generate a complete 2D segmentation, we fill these holes in using the SLIC algorithm, resulting in a complete segmentation, seen in the top right. The bottom right shows human annotated ground truth for the scene.



**Figure 3.6.3:** Examples of under-segmentation output. From left to right- ground truth annotation, SLIC, GCb10, DASP, and VCCS. Each is shown with two different superpixel densities.



**Figure 3.6.4:** Boundary recall and under-segmentation error for SLIC, GCb10, DASP, and VCCS.

are shown in Figure 3.6.2.

### 3.6.2 SUPERVOXELS

In order to evaluate the quality of supervoxels generated by VCCS, we performed a quantitative comparison with three state-of-the-art superpixel methods using publicly available source code. We selected the two 2D techniques with the highest published performance from a recent review [5]: a graph based method, GCb10 [84]<sup>3</sup>, and a gradient ascent local clustering method, SLIC [5]<sup>4</sup>. Additionally, we selected another method which uses depth images, DASP[89]<sup>5</sup>. Examples of over-segmentations produced by the methods are given in Figure 3.6.3.

#### OBJECT BOUNDARY ADHERENCE

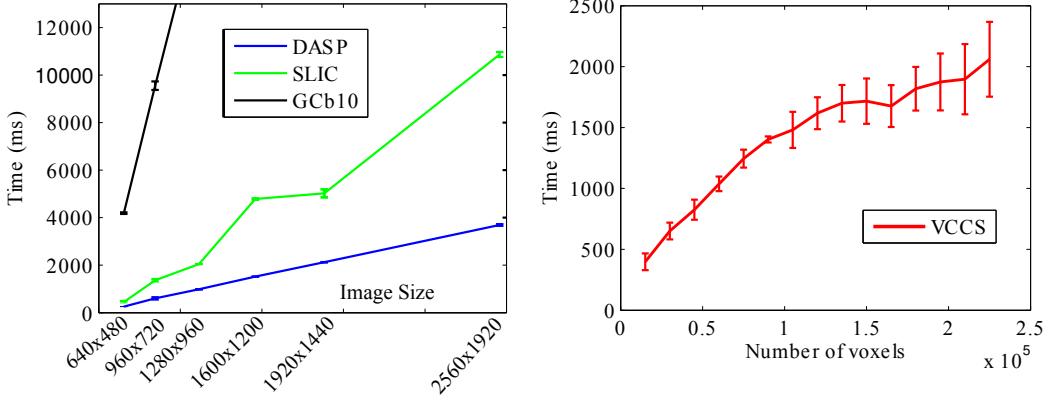
The most important property for superpixels is the ability to adhere to, and not cross, object boundaries. To measure this quantitatively, we have used two standard metrics for boundary adherence- boundary recall and under-segmentation error[54, 84]. Boundary recall measures what fraction of the ground truth edges fall within at least two pixels of a superpixel boundary. High boundary recall indicates that the superpixels properly follow the edges of objects in the ground truth labeling. The results for boundary recall are given in Figure 3.6.4. As can be seen, VCCS and SLIC have the best boundary recall performance, giving similar results as the number of superpixels in the segmentation varies.

Under-segmentation error measures the amount of leakage across object boundaries. For a ground truth segmentation with regions  $g_1, \dots, g_M$ , and the set of superpixels from an over-

<sup>3</sup><http://www.csd.uwo.ca/~olga/Projects/superpixels.html>

<sup>4</sup>[http://ivrg.epfl.ch/supplementary\\_material/RK\\_SLICSuperpixels/index.html](http://ivrg.epfl.ch/supplementary_material/RK_SLICSuperpixels/index.html)

<sup>5</sup><https://github.com/Danvil/dasp>



**Figure 3.6.5:** Speed of segmentation for increasing image size and number of voxels. Use of GCb10 rapidly becomes unfeasible for larger image sizes, and so we do not adjust the axes to show its run-time. The variation seen in VCCS run-time is due to dependence on other factors, such as  $R_{seed}$  and overall amount of connectivity in the adjacency graphs.

segmentation,  $s_1, \dots, s_K$ , under-segmentation error is defined as

$$E_{useg} = \frac{1}{N} \left[ \sum_{i=1}^M \left( \sum_{s_j | s_j \cap g_i} |s_j| \right) - N \right], \quad (3.9)$$

where  $s_j | s_j \cap g_i$  is the set of superpixels required to cover a ground truth label  $g_i$ , and  $N$  is the number of labeled ground truth pixels. A lower value means that less superpixels violated ground truth borders by crossing over them. Figure 3.6.4 compares the four algorithms, giving under-segmentation error for increasing superpixel counts. VCCS outperforms existing methods for all superpixel densities.

## TIME PERFORMANCE

As superpixels are used as a preprocessing step to reduce the complexity of segmentation, they should be computationally efficient so that they do not negatively impact overall performance. To quantify segmentation speed, we measured the time required for the methods on images of increasing size (for the 2D methods) and increasing number of voxels (for VCCS). All measurements were recorded on an Intel Core i7 3.2Ghz processor, and are shown in Figure 3.6.5. VCCS shows performance competitive with SLIC and DASP (the two fastest superpixel methods in the literature) for voxel clouds of sizes which are typical for Kinect data at  $R_{voxel} = 0.008m$  (20-40k voxels). It should be noted that only VCCS takes advantage of multi-threading (for octree, kd-tree, and FPFH computation), as there are no publicly available multi-threaded implementations of the other algorithms.



**Figure 3.6.6:** Example results for the OSD dataset. Points beyond a distance of 2m were cropped for visualization. Parameters:  $R_{voxel} = 0.005$ ,  $R_{seed} = 0.02$ ,  $\beta_{\text{Thresh}} = 10^\circ$ .

### 3.6.3 LOCALLY CONVEX CONNECTED PATCHES

We compare segments found using LCCP against ground truth using three standard measures: *Weighted Overlap* (*WOv*), which is a summary measure proposed by Silberman *et al.* [79], as well as *false negative* (*fn*) and *false positive* (*fp*) scores from [82] and *over-* ( $F_{os}$ ) and *under-segmentation* ( $F_{us}$ ) from [72].

Method	Learned Features	WOv		<i>fp</i>		<i>fn</i>		$F_{os}$	$F_{us}$
		Mean	Mean	SD	Mean	SD	Mean	Mean	Mean
LCCP	NO LEARNING	88.7%	4.8%	2.6%	8.3%	8.7%	7.4%	4.7%	
Richtsfeld [72]	RGB-D,Texture,Geometry	-	-	-	-	-	4.5%	7.9%	
Ückermann [82]	NO LEARNING	-	1.9%	3.3%	7.8%	7.3%	-	-	

**Table 3.6.1:** Comparison of different segmentation methods on the OSD dataset using weighted overlap *WOv* (the higher, the better), false positives  $f_p$ , false negatives  $f_n$ , as well as over- and under-segmentation  $F_{os}$  and  $F_{us}$  (the lower, the better). LCCP results were produced with voxel resolution  $R_{voxel} = 0.005$ , seed resolution  $R_{seed} = 0.02$  and concavity tolerance angle  $\beta_{\text{Thresh}} = 10^\circ$ .

The qualitative examples from the OSD dataset (Figure 3.6.6) show that LCCP performs very well in the segmentation of cluttered scenes. The object separation can be intuitively understood: all objects present in the scenes are separated by concave boundaries, i.e. a line connecting neighboring surfaces of two different objects always travels through “air”. This is also true for the boundary between an object and the supporting surface. As a consequence, objects that have a convex shape are correctly captured as one segment and separated from the other objects. Hollow objects (bowls, cups etc.) can be observed to show multiple segments inside, because the orientation of surface normals changes strongly on these concave surfaces.

Method	Learned Features	Depth Data	WOv
LCCP	NO LEARNING	depth	53.6%
	NO LEARNING	smoothdepth	53.8%
LCCP + ext. convexity	NO LEARNING	smoothdepth	57.6%
Silberman <i>et al.</i> [79]	RGB	-	50.3%
	Depth	both	53.7%
	RGB-D	both	60.1%
	RGB-D + Support + Structure classes	both	61.1%
Gupta <i>et al.</i> [37]	gPb-ucm Gradients (from [11])	-	55.0%
	gPb-ucm + Depth + Concavity Gradients	both	62.0%

**Table 3.6.2:** Comparison of different segmentation methods on the NYU dataset using weighted overlap  $WOv$ . LCCP results were produced with voxel size  $R_{voxel} = 0.01$ , seed size  $R_{seed} = 0.04$  and concavity tolerance angle  $\beta_{\text{thresh}} = 10^\circ$ .

The quantitative results (Table 3.6.1) demonstrate that the approach is able to compete with state-of-the-art methods in the task of segmenting cluttered scenes with ‘single-part’ objects.

Example scenes in Figure 3.6.1 show that the LCCP also works well on the real-world scenes from the NYU dataset. The quantitative results (Table 3.6.2)<sup>6</sup> show that our algorithm is able to produce good results on the challenging dataset. Despite being much simpler and without requiring learning on human annotated ground-truth, we compete with the approach from [79] when only depth information is used. Additionally, we still achieve 93% of their score when comparing against the more complex feature spaces used in conjunction with learning-based algorithms. We should emphasize that our competitors do not aim for object parts but rather for “whole object” detections, specifically, those whole objects learned from this particular annotated ground truth. Conversely, our method establishes a general rule for object-ness that does not depend on this particular dataset, nor on the whims of a particular human annotator.

### 3.7 DISCUSSION

In this Chapter we have presented several new concepts- the octree adjacency graph, supervoxels, a segmentation method which uses supervoxels, as well as a way to sequentially update an octree with new frames of data. Additionally, we have presented quantitative and qualitative results which demonstrate the usefulness of these techniques on real-world datasets. In particular, LCCP has demonstrated the usefulness of a patch-based adjacency-graph interpretation of 3D Point Cloud data. The results we achieved stem from two core properties: the ability of supervoxels to efficiently encode local regions, and the usefulness of a 3D adjacency-graph in resolving situations which are ambiguous in a 2D representation. Additionally, the

<sup>6</sup>Updated results for [79] are available at [http://cs.nyu.edu/~silberman/datasets/nyu\\_depth\\_v2.html](http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html).

depth-dependent octree world model is able to compactly represent 3D data with included adjacency information that gracefully adjusts level-of-detail as distance increases from the sensor. Finally, the sequential update process we described allows a basic “object permanence” to be directly encoded without the need for higher-level objects. While this does have the advantage of giving “voxel-permanency”, it does suffer from an inability to deal with the case of “moving-while-occluded”, a situation which cannot be resolved at the low-level used here.

In the next Chapter, we will use the octree model we have developed here as observations for particle filter tracking. We will demonstrate its effectiveness as the basis for correspondence association, and show how its preservation of voxels through occlusions allows us to track objects in a real-world application. Additionally, we shall use the supervoxels we presented here as the basis for dividing tracked models into strata, and show how this can be used to vastly improve the run-time performance of a tracker.

*Optical Illusion is optical truth! ... In them is evidenced the living interaction of our inner nature with outer nature.*

Johann Wolfgang von Goethe

# 4

## Model-Based Point Cloud Tracking

NOW THAT WE HAVE ESTABLISHED a reduced, stable world model in which voxels persist through occlusions, the next step is to adapt the general framework of Sequential Bayesian Filtering (SBF) to track models within this 3D voxel world. For an introduction to the general framework of SBF, we refer the reader to Appendix B). In this Chapter we begin by presenting the basic framework of particle filter tracking in 3D point clouds, and show how point correspondences can be used to evaluate the particle filtering likelihood function.

While the correspondence approach is feasible for tracking single targets, it suffers from the same flaw as other approaches when extended to multiple targets [46, 48, 52] - it significantly increases the computational resources required. This increase is due to the need for more particles - due to assignment of particles to individual targets, a larger state space, or independent filters for each target. While there has been work addressing this problem by offloading processing to a GPU [26], in this work we take a different approach, and search for fundamental changes to the point cloud correspondence particle filter which can reduce computational complexity without affecting accuracy.

This Chapter is organized as follows: First, in Section 4.1 we present a framework for particle filter tracking in point clouds using point correspondences. In Section 4.2 we present the primary contribution of this Chapter; the use of a supervoxel-based stratified sampling approach to greatly reduce the computational complexity of point cloud correspondence particle filtering. Finally, in Section 4.3 we will show that the approach allows performance (on a

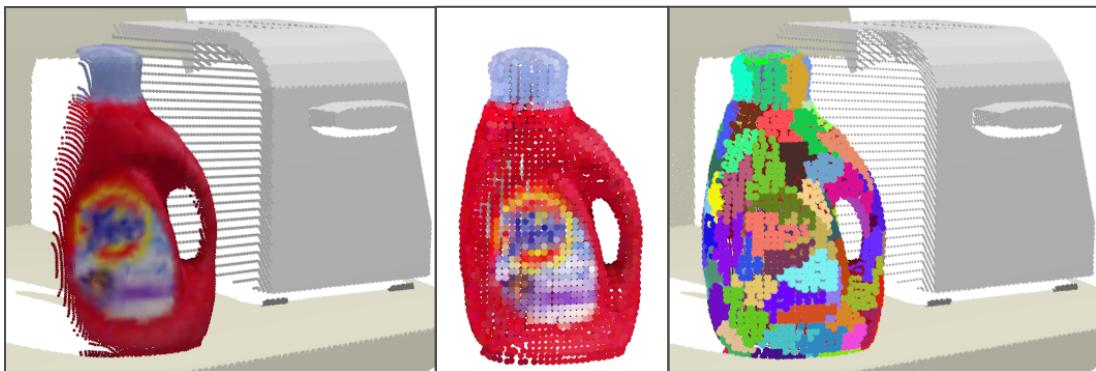
standard CPU) exceeding that which can be obtained on a recent GPU implementation [26]. Furthermore, we shall present extensive quantitative experiments demonstrating the benefits of this approach, as well as show qualitative results from a real-world application.

## 4.1 PARTICLE FILTERS IN 3D

The underlying mechanics of 3D point cloud correspondence particle filtering remain the same as in Chapter 2, and so we shall not discuss them extensively here; for a detailed introduction to the topic, we again refer the reader to [30] or [88]. Rather, we shall only discuss the aspects that differentiate it - the models and how they are scored and propagated. The models here consist of point clouds, and the measurement function relies on point to point correspondence for scoring, rather than a global per-detection metric (such as a histogram distance, used in the 2D trackers of in Chapter 2). The dynamic model uses real-world 3D coordinates which also include orientation, rather than 2D pixel coordinates in the image plane. The primary novelty of the approach we present here lies in how we score individual particle predictions using the measurement model.

### 4.1.1 MODEL REPRESENTATION

One of the main limitations of the 2D projected mask model discussed in Chapter 2 is that the masks of objects are not invariant to pose changes - in general, rotation of an object will change the shape of its mask and distribution of its color histogram. As we now have the ability to observe the full 3D shape of an object, we choose to represent objects as clusters of points which correspond to the exterior of the object. A visual representation of such a model is given in Figure 4.1.1.



**Figure 4.1.1:** Example of data from “Tide” sequence. The left frame shows an example of the raw input cloud. Sampling effects from the synthetic RGB-D camera are visible in the quantization of points, especially on the edges of objects. The middle frame shows the voxelized model representation we use, while the right frame shows an example of supervoxel strata used for sampling with  $R_{seed} = 0.07m$ .

Points for objects are stored in a model-centered reference frame (which we shall denote with superscript  $m$ ), with each point containing an XYZ position, an RGB color for the point, as well as a surface normal vector. That is, each point  $p$  of the model  $k$  consists of a nine-dimensional vector:

$$p_k^m = [x^m, y^m, z^m, R, G, B, n_x, n_y, n_z], \quad (4.1)$$

and a model for an object  $O_k$  consists of a vector of  $n_k$  such points  $p^m$ :

$$O_k^m = [p_0^m \dots p_{n_k}^m]. \quad (4.2)$$

It is important to note that the points of an object model given above are model-relative - they must be transformed into the world coordinates in order to evaluate their fit to observations. This will be discussed further in the next Section.

#### 4.1.2 DYNAMIC MODEL

In the 2D tracker presented previously, the time-dependent state vector of a particle consisted of a position shift vector  $\mathbf{p}_t = [p_x, p_y]$  and a velocity vector  $\mathbf{v}_t = [v_x, v_y]$ . The natural extension of this to 3D is to simply add a third  $p_z$  and  $v_z$  element to each. Of course we should note that the  $x$  and  $y$  dimensions here in our 3D representation are distinct from those in 2D, which represented pixel coordinates in the image plane. Here our positional coordinates represent real-world distances from a fixed origin (typically the camera “pin-hole” position). It is also important to note that coordinates in our 3D representation are originally in a continuous space - though we discretize them using the octree model discussed in the previous Chapter. For clarity, we shall simply denote coordinates in the world reference frame with no superscript.

While this straightforward extension gives us a reasonable 3D equivalent to our 2D tracked masks, we now have full 3D models, and so it makes sense to use a state vector which takes advantage of it. As such, we further extend the state vector for position and velocity to allow for rotations of the model around the object reference frame x-axis (roll -  $\gamma$ ), y-axis (pitch -  $\beta$ ), and z-axis (yaw -  $\alpha$ ). This yields a position state vector for particle  $j$  at time  $t$  of

$$\mathbf{x}_t^j = [d_x, d_y, d_z, \gamma, \beta, \alpha]. \quad (4.3)$$

Each object model is tracking using a set of  $N$  such particles. We shall now generally omit the object variable  $k$  in our notation for clarity. Even though we omit the  $k$ , the reader should assume that the following equations are for individual object models, and that we have a set of  $N$  independent particles for each object. Additionally, we have velocity state vector

$$\mathbf{v}_t = [v_x, v_y, v_z, v_\gamma, v_\beta, v_\alpha], \quad (4.4)$$

which is not tracked individually per particle, but rather as a whole for the model.

As before, motion is modeled using a constant velocity model in discrete time with a variable sampling period  $T$ , giving the dynamic model

$$\mathbf{x}_t = \mathbf{x}_{t-1} + T\mathbf{v}_{t-1} + \boldsymbol{\omega}, \quad (4.5)$$

with noise vector  $\boldsymbol{\omega}$  assumed to be zero mean Gaussian with fixed covariance. Particle velocities are updated after weighting of individual particles using the measurement model, and are a weighted average of the change in position

$$\mathbf{v}_t = \frac{1}{TN} \sum_{j=1}^N w_j (\mathbf{x}_t^j - \mathbf{x}_{t-1}^j), \quad (4.6)$$

where  $w_j$  is the normalized weight for particle  $j$ .

Tracking independent velocities for each particle doubles the dimensionality of the state-space, requiring a proportional increase in the number of particles. While the use of independent velocity states potentially helps in complicated tracking scenarios, in our experiments we were unable to observe any tangible benefit. Moreover, in order to avoid instability in the tracking results we needed to double the number of particles for a given noise level, doubling the processing time required. As such, we have chosen to use the above “group-velocity”, and leave it to future work to investigate the possibility of independent velocity states.

#### 4.1.3 MEASUREMENT MODEL

As points for the model are given in a model-centered frame of reference, we must transform them to the world frame them using a 3D affine transformation quaternion:

$$\mathbf{B}^j = \begin{bmatrix} \cos \alpha \cos \beta & \cos \alpha \sin \beta \sin \gamma - \sin \alpha \cos \gamma & \cos \alpha \sin \beta \cos \gamma + \sin \alpha \sin \gamma & d_x \\ \sin \alpha \cos \beta & \sin \alpha \sin \beta \sin \gamma + \cos \alpha \cos \gamma & \sin \alpha \sin \beta \cos \gamma - \cos \alpha \sin \gamma & d_y \\ -\sin \beta & \cos \beta \sin \gamma & \cos \beta \cos \gamma & d_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (4.7)$$

which we use to transform the extended position vector for each point in the model:

$$p^m = [x^m, y^m, z^m, 1], \quad (4.8)$$

yielding positions in the world frame for each of our  $\eta$  model points for a particular particle  $j$ :

$$\begin{bmatrix} \mathbf{p}_1^j \\ \mathbf{p}_2^j \\ \vdots \\ \mathbf{p}_\eta^j \end{bmatrix} \begin{bmatrix} [x_1, y_1, z_1, 1]^\top \\ [x_2, y_2, z_2, 1]^\top \\ \vdots \\ [x_\eta, y_\eta, z_\eta, 1]^\top \end{bmatrix} = \begin{bmatrix} \mathbf{B}^j & 0 & \dots & 0 \\ 0 & \mathbf{B}^j & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{B}^j \end{bmatrix} \begin{bmatrix} [x_1^m, y_1^m, z_1^m, 1]^\top \\ [x_2^m, y_2^m, z_2^m, 1]^\top \\ \vdots \\ [x_\eta^m, y_\eta^m, z_\eta^m, 1]^\top \end{bmatrix}. \quad (4.9)$$

Once we have our transformed points, we then must establish correspondences between each particle's model points and a world point. This is done so that we may score how well a particular particle matches the current world model observation. That is, for each transformed point  $\mathbf{p}_{1\dots\eta}^j$ , we select corresponding point  $\mathbf{p}^*$  in the observation which has minimal spatial distance. To find these correspondences, we first compute a KD-tree in the spatial dimensions for the world model points. This allows us to efficiently search for the nearest point to each transformed point. We create this tree for the world model rather than the transformed model (even though the former has more points) as there is only one world, but many particles and models. Computing it for the models would require a KD-tree for each particle in each model. Additionally, computing it for the world allows us to take advantage of sampling strategies (discussed in the next Section) which significantly reduce our run-time complexity.

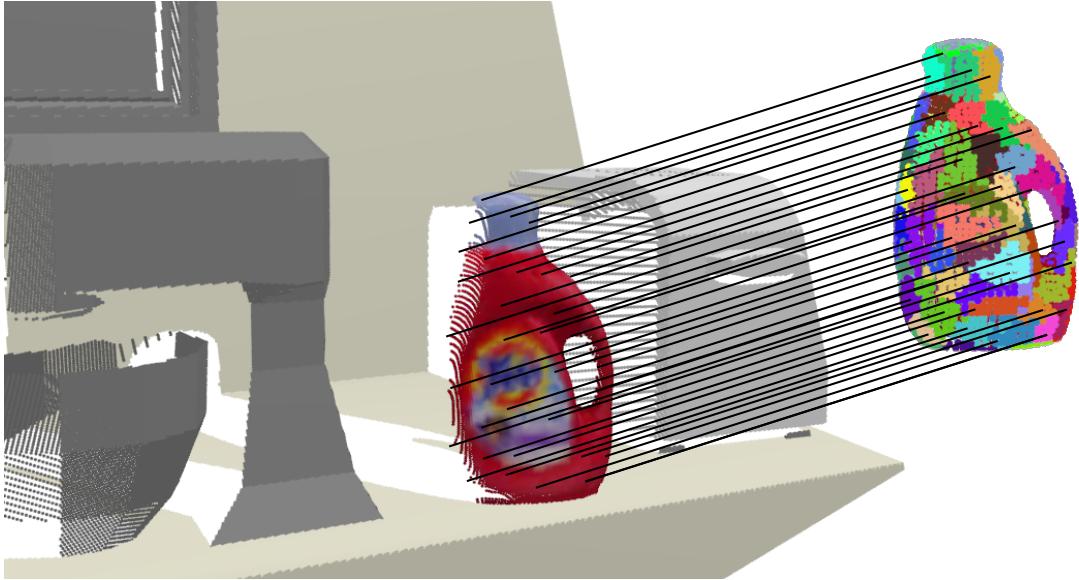
Once we have selected (with replacement) an observed point correspondence for each model point, we must calculate an un-normalized weight  $\tilde{w}^j$  corresponding to the similarity of the transformed points to the world observation. This is accomplished by summing the individual correspondence scores computed using weighted distance in world-, color-, and normal-space:

$$\tilde{w}^j = \sum_1^\eta \frac{1}{1 + \frac{\mu \|\mathbf{p}_{xyz}^j - \mathbf{p}_{xyz}^*\|}{R_{voxel}} + \frac{\lambda D_c(p_{RGB}^j, p_{RGB}^*)}{m} + \varepsilon \|\mathbf{p}_{n_x n_y n_z}^j - \mathbf{p}_{n_x n_y n_z}^*\|}, \quad (4.10)$$

where we follow the convention given in Section 3.2.2. That is,  $\mu$ ,  $\lambda$ , and  $\varepsilon$  are weighting constants,  $D_c$  is euclidean distance in HSV space, and  $m$  is a normalizing constant. We do not normalize normals, as they are already unit vectors. In our experiments we typically set the weighting factors to  $\mu = 1$ ,  $\lambda = 2$ ,  $\varepsilon = 1$ , as this balances the scoring between color and geometric shape, and found experimentally that it produced consistently good tracking results. The calculated particle weights  $\tilde{w}^j$  are then normalized, and a final state estimate can be computed by taking the weighted average of all particles

$$\mathbf{x}_t = \sum_{j=1}^N w_j \mathbf{x}_t^j, \quad (4.11)$$

and the group-velocity can be computed using Equation 4.6.



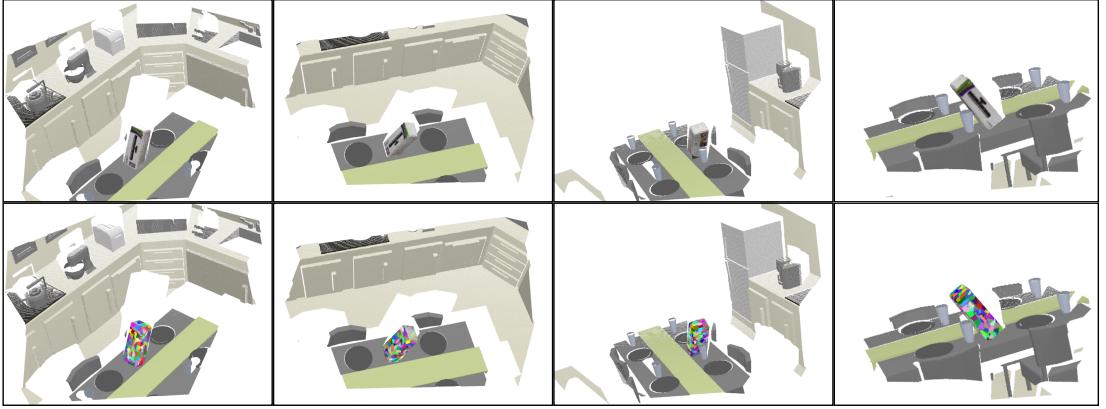
**Figure 4.2.1:** The model is divided into strata (shown as separate colors), and a random point is selected from each stratum for correspondence matching. The colored circles overlaid on the observed point cloud show the search radius used for finding correspondences.

## 4.2 STRATIFIED CORRESPONDENCE SAMPLING

While the tracking methodology discussed above works, in practice its run time performance is very poor, even for single objects. Moreover, speed of tracking is highly dependent on the size of object models as well as voxel resolution used. To address this, we propose a sampling scheme which selects a limited number of points from the model to transform and test. By doing this, we achieve linear asymptotic time complexity for the particle filter with respect to the number of particles - there is no dependence on the number of points in the models or the voxel resolution used. The only step which is dependent on the number of input points is the KD-tree construction, but this is only done once for the world model independent of the number of trackers, and must be done as a pre-processing step anyway for normal computation.

The proposed sampling scheme is as follows. We select a spatial sampling resolution  $R_{seed}$  based on the number of desired sample points per particle  $N_s$ . We then divide the model into strata, where each stratum is a supervoxel using the VCCS method described previously [62]. Supervoxels are a voxel-based surface patch representation that use connectivity, colors, and normals so that their edges conform well to object part boundaries. The strata are evenly divided over the spatial structure of the model, as seen in Figure 4.1.1. Additionally, using supervoxels as the strata ensures that we sample the important features of the models - for example in the model of Figure 4.1.1, we have a stratum for the brand logo, as well as ones for the concavities of the handle.

For each particle, we randomly select a point from each stratum using uniform sampling,



**Figure 4.2.2:** Tracking on the artificial “Kinect Box” sequence. The top row shows tracked output overlaid on input data, while the bottom row shows the supervoxel strata that are used for sampling.

and then transform and score it as described in the previous Section. As an additional step, we also select  $\frac{N_s}{4}$  points uniformly from the entire model. Using strata reduces the noise which occurs when sampling from the whole model exclusively, while sampling randomly from the entire distribution improves occlusion performance.

While sampling will tend to produce noisier tracking results for low  $N_s$ , it also greatly reduces the computational complexity, as we only need to transform and test a small subset of the model points. This allows one to greatly increase the number of particles for a given frame-rate. Importantly, each particle is testing a separate random subset of model points. This results in the product of  $N_s$ , the number of sample points per particle, and  $N$ , the number of particles, reaching a critical level where coverage becomes sufficient that error is equivalent to sampling all points. In the results presented below, we shall demonstrate that this critical level can be used to significantly decrease run time for a given level of error. That is, we shall show that the number of points that must be tested overall, for a given level of error, is lower when stratified sampling is used. This means that we can significantly increase accuracy for a given frame-rate, reducing run-time complexity to the point that we can track 6 DoF pose for multiple objects in real-time.

### 4.3 EXPERIMENTAL RESULTS

In this Section we first present results on a set of synthetic videos to quantify the effect of the stratified sampling, and compare results to a state of the art GPU particle filter [26]. We then present qualitative results on real videos in a robotic learning application, where we track multiple interacting targets with significant occlusions. In both synthetic and real cases, input consists of RGB-D sequences. Trackers were initialized using an external pose - in the synthetic case, from ground truth, and in the real case, using a pose estimation algorithm [22]. Object models were generated by registering multiple views of the objects using the same RGB-D sen-

sor employed for tracking. All experiments were performed on a standard desktop computer (Intel i7 3.2Ghz), using all four available cores.

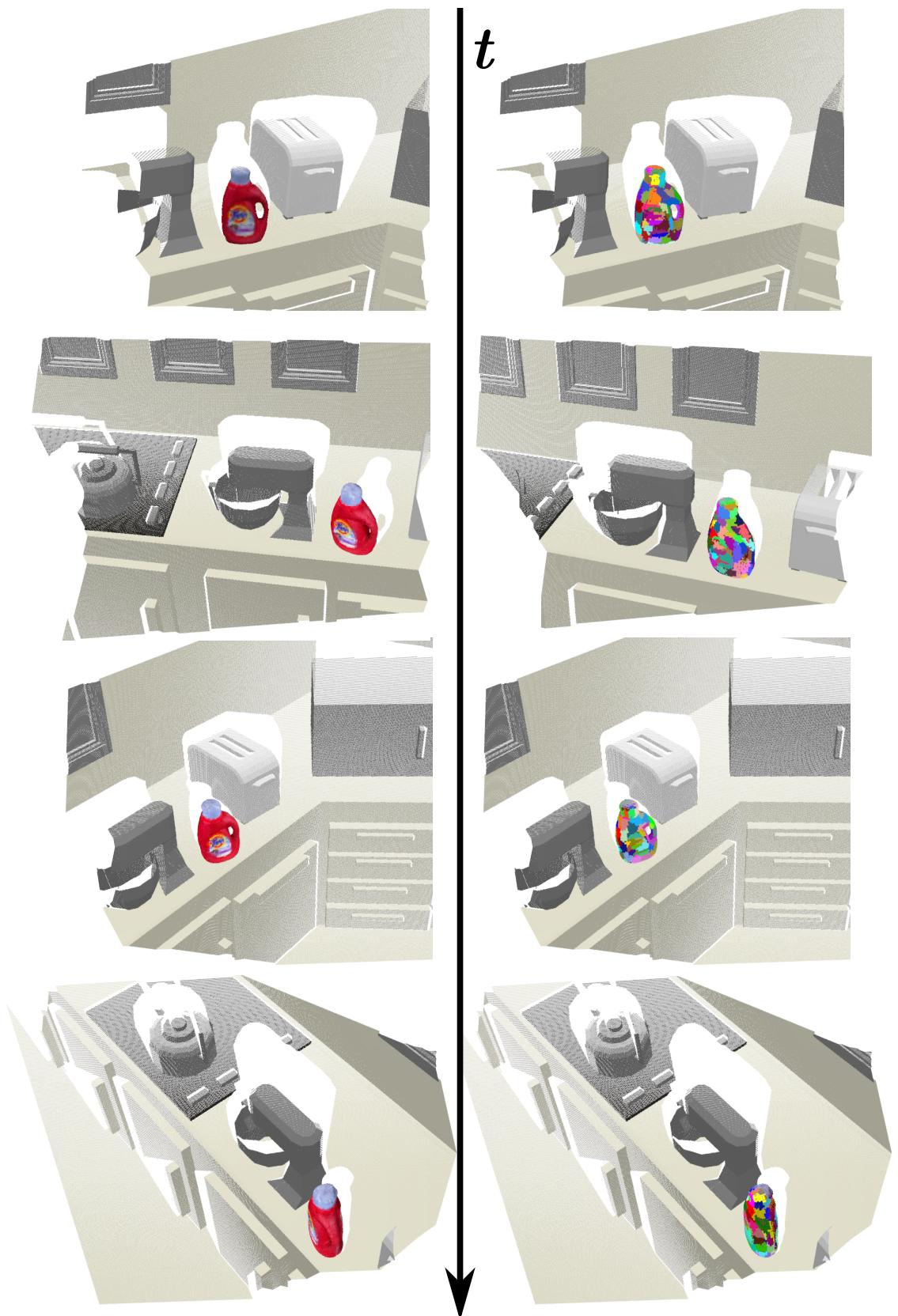
#### 4.3.1 RESULTS ON SYNTHETIC SEQUENCES

In our first experiment, we demonstrate the effectiveness of our stratified sampling strategy using four synthetic tracking videos from [26]. These RGB-D sequences are set in a virtual kitchen (see Figure 4.2.2) and each contain a single item to track as the camera moves. Ground truth trajectories of the cameras are given in Figures 4.3.2-4.3.8; one can observe that the trajectories are complex, consisting of large variations in position, orientation, and velocity.

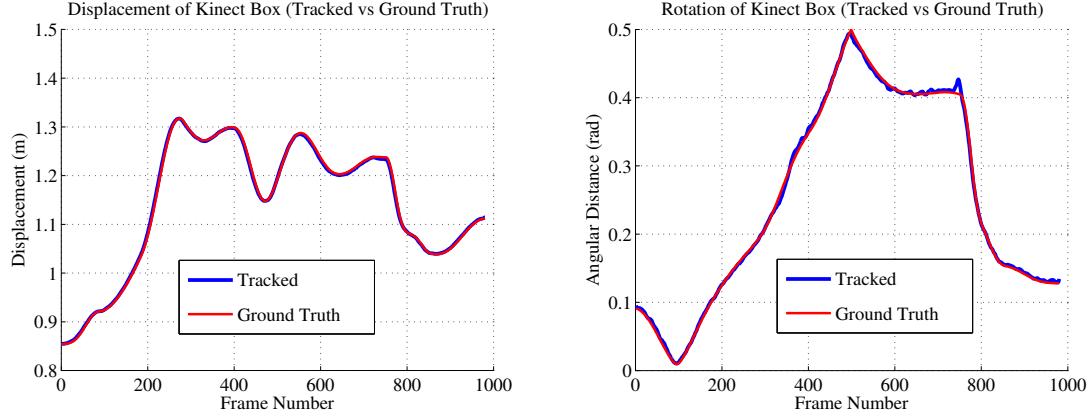
To evaluate our approach, we compute root mean square (RMS) error in both translation and orientation, averaged over 25 test runs for each sequence. Computation times are measured in ms per frame, and are also averaged across all frames of the 25 test runs. In order to compare with [26], we have combined their RMS error results for each dimension (x, y, z, roll, pitch, yaw) into two measurements - displacement and rotation. Rotation is calculated using the unit quaternion distance metric [49], which is equivalent to the angular distance on the unit sphere. This combination reduces the amount of data to compare without loss, as the choice of orientation of the dimensions is arbitrary and without import. Example displacement and rotation ground truths for the “Kinect Box” sequence can be found in Figure 4.3.2.

Timing results are given in Figures 4.3.3 -4.3.9, showing results for the four sequences, with each plot scanning across number of particles and number of sample points. Results for all four sequences are similar. One can observe that, for a given level of sampling, the RMS error decreases for both displacement and rotation as the number of particles increases. More importantly, it is also apparent that, for a given level of error, run-time per frame can be minimized by reducing the number of samples used and increasing the number of particles. Additionally, one can observe that RMS error appears to be asymptotic, with lower sampling levels approaching the asymptote at lower run-times.

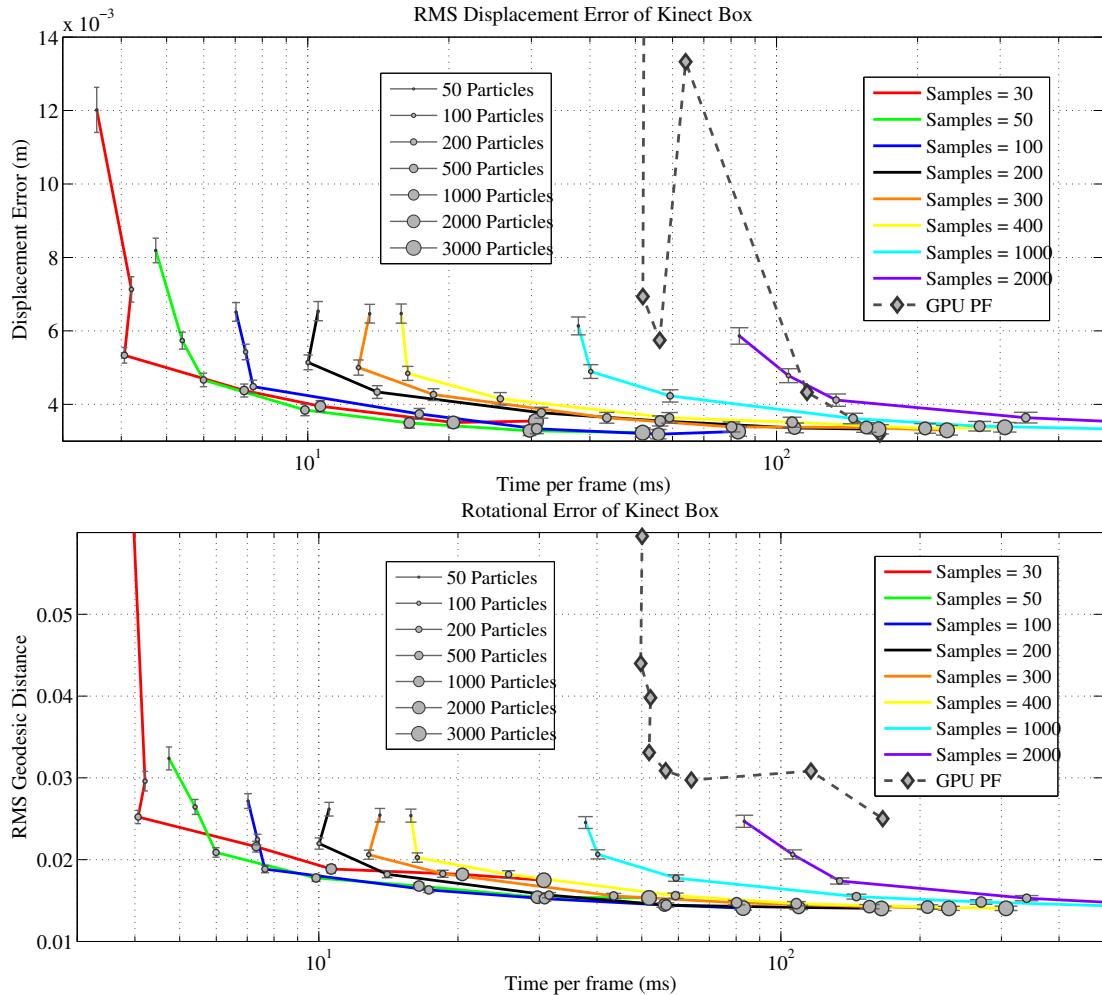
We should also note that the minimum error asymptote observed is likely a consequence of the sampling resolution of the synthetic Kinect camera. For example, in the “Kinect Box” sequence, average distance to neighboring points (8-neighborhood) on the tracked box surface is 3.3 mm. This corresponds almost exactly to our observed error asymptote. This can be observed in all four sequences - our minimal error corresponds closely to the average point to point resolution of the observations on the model.



**Figure 4.3.1:** Tracking on the artificial “Tide” sequence. The left column shows tracked output overlaid on input data, while the right column shows the supervoxel strata that are used for sampling.

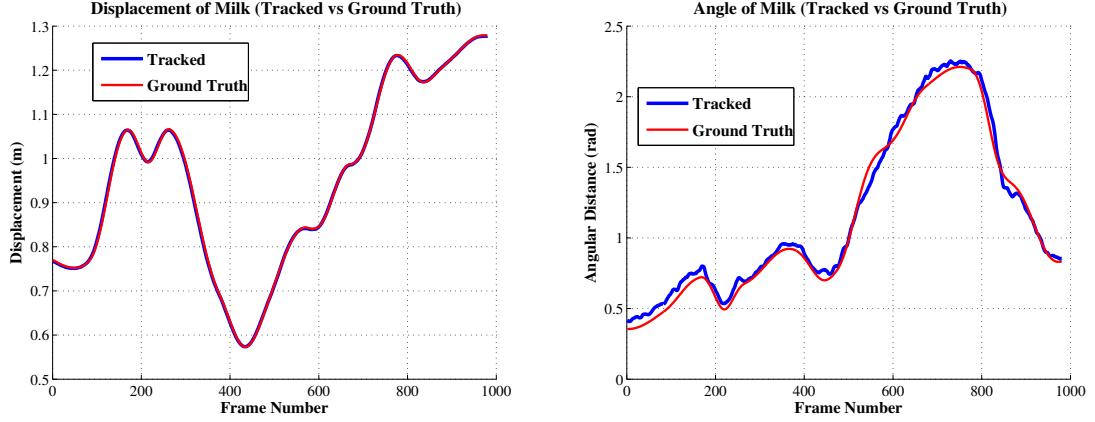


**Figure 4.3.2:** Displacement and rotation ground truth, with an example tracked result from a single run at  $N_{samples} = 100$  and  $N_{particles} = 1000$  (a frame rate of 20 fps).

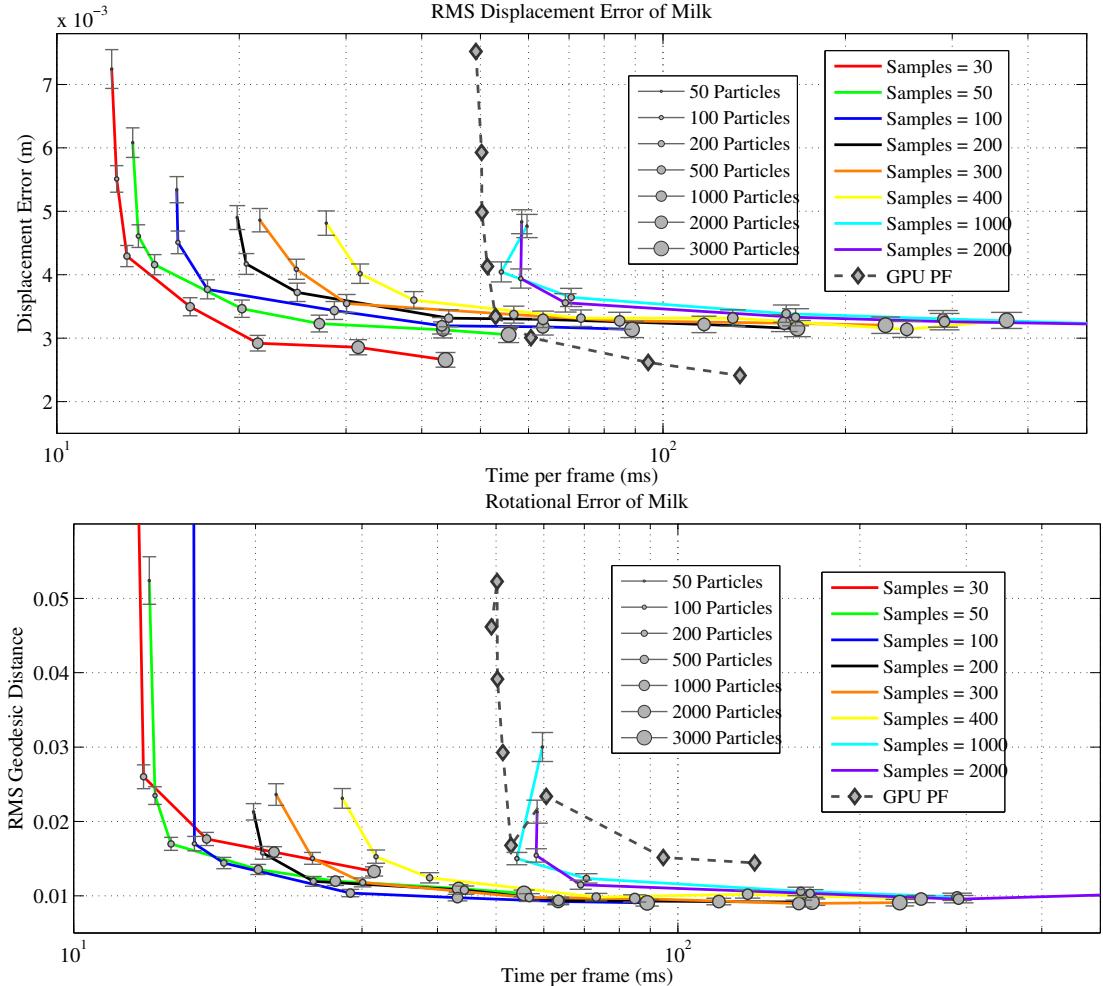


**Figure 4.3.3:** Results on the Kinect Box artificial sequence. Each colored curve represents a certain number of samples, and gives mean RMS error averaged over 25 trial runs for increasing numbers of particles.

Our performance compares favorably to the results of Choi and Christensen [26] - for a given level of error, we achieve per-frame run times that are between half and a tenth of their published results. Additionally, we consistently reach the error asymptote at considerably

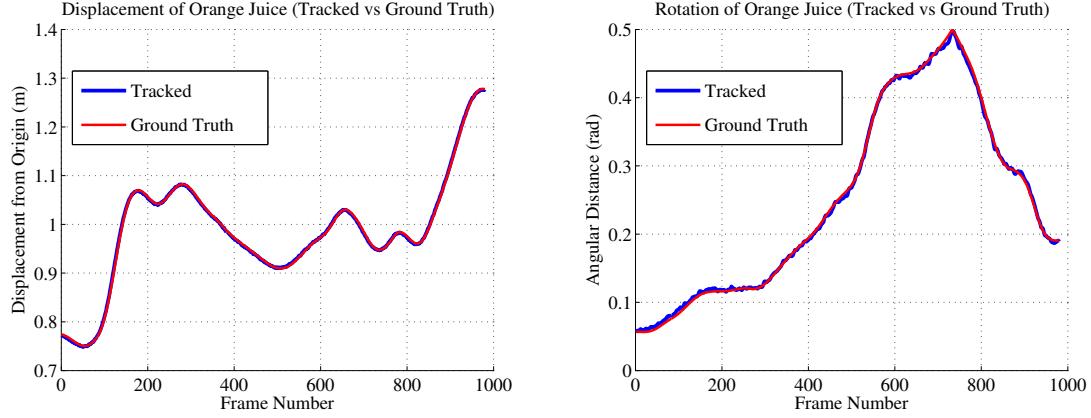


**Figure 4.3.4:** Displacement and rotation ground truth, with an example tracked result from a single run at  $N_{samples} = 100$  and  $N_{particles} = 1000$  (a frame rate of 20 fps).

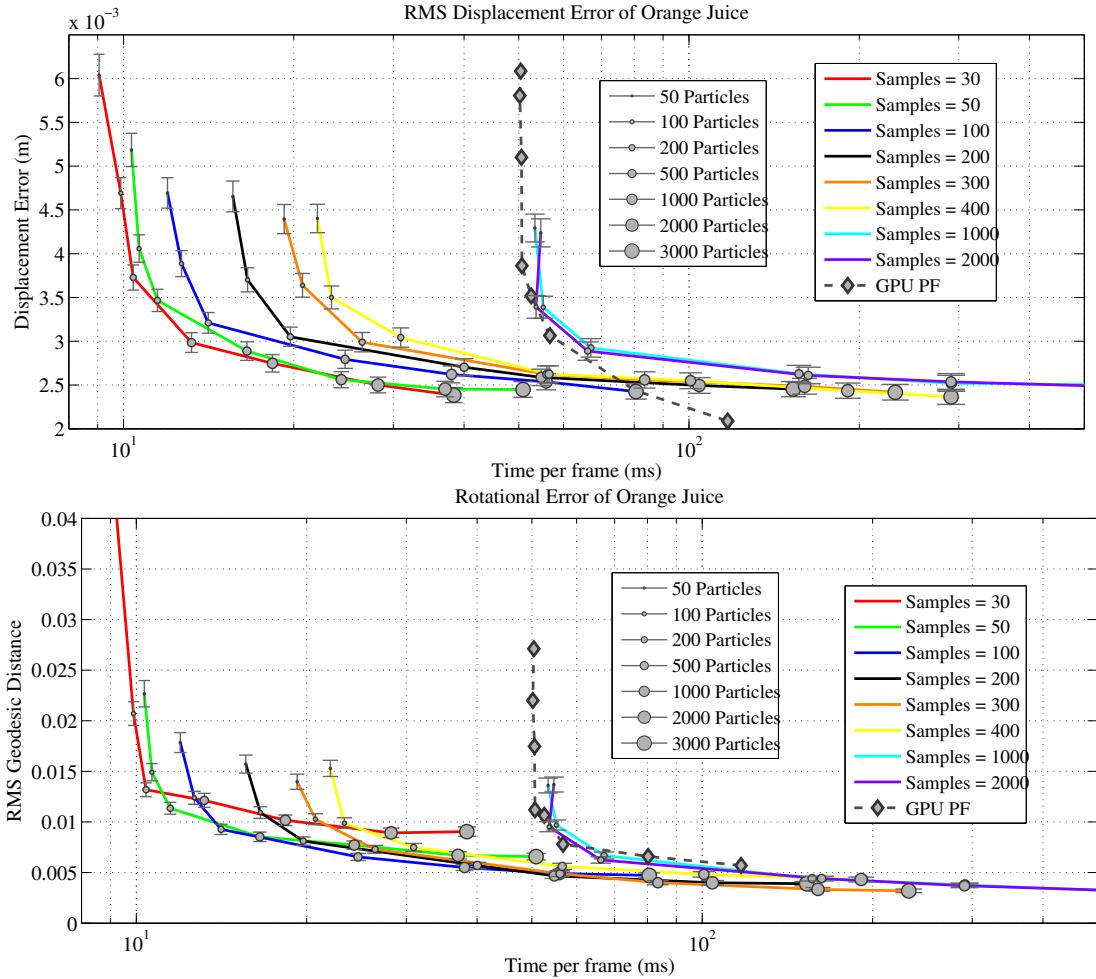


**Figure 4.3.5:** Results on the Milk artificial sequence. Each colored curve represents a certain number of samples, and gives mean RMS error averaged over 25 trial runs for increasing numbers of particles.

lower run times. We should also note that the highest sampling level shown corresponds to a complete sampling of the model, and can be thought of as equivalent to the baseline PCL implementation, although we have made some modifications to the resampling and dynamic

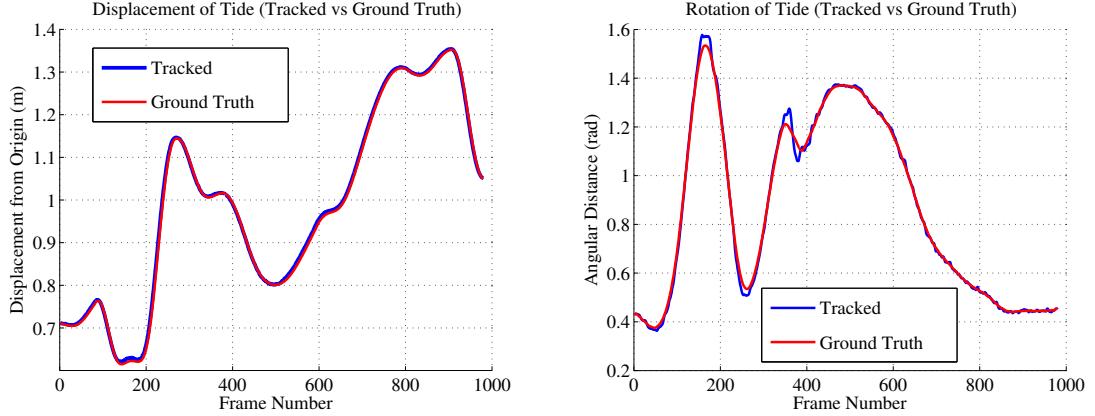


**Figure 4.3.6:** Displacement and rotation ground truth, with an example tracked result from a single run at  $N_{samples} = 100$  and  $N_{particles} = 1000$  (a frame rate of 20 fps).

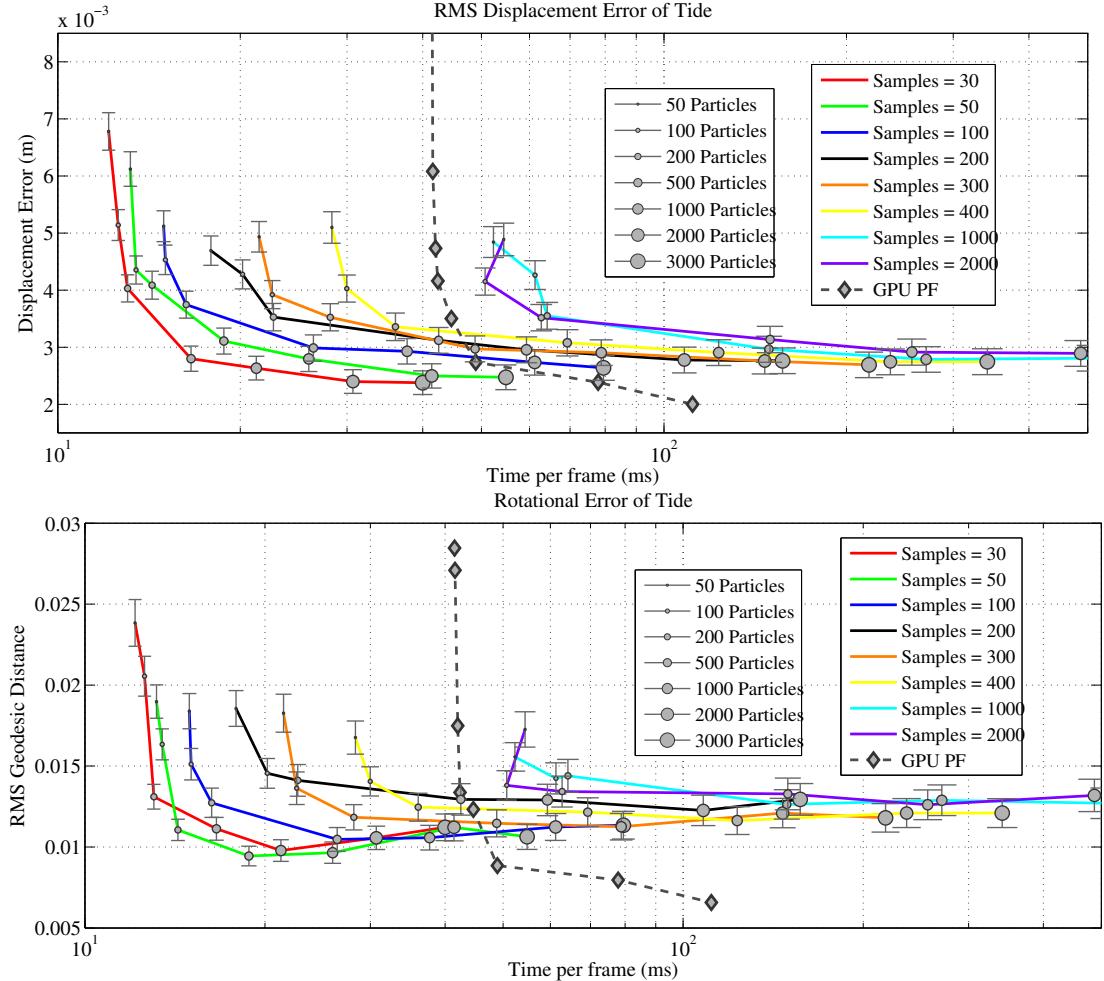


**Figure 4.3.7:** Results on the Orange Juice artificial sequence. Each colored curve represents a certain number of samples, and gives mean RMS error averaged over 25 trial runs for increasing numbers of particles.

model which improve results. As can be seen, we are at least an order of magnitude faster than this base implementation.



**Figure 4.3.8:** Displacement and rotation ground truth, with an example tracked result from a single run at  $N_{samples} = 100$  and  $N_{particles} = 1000$  (a frame rate of 20 fps).

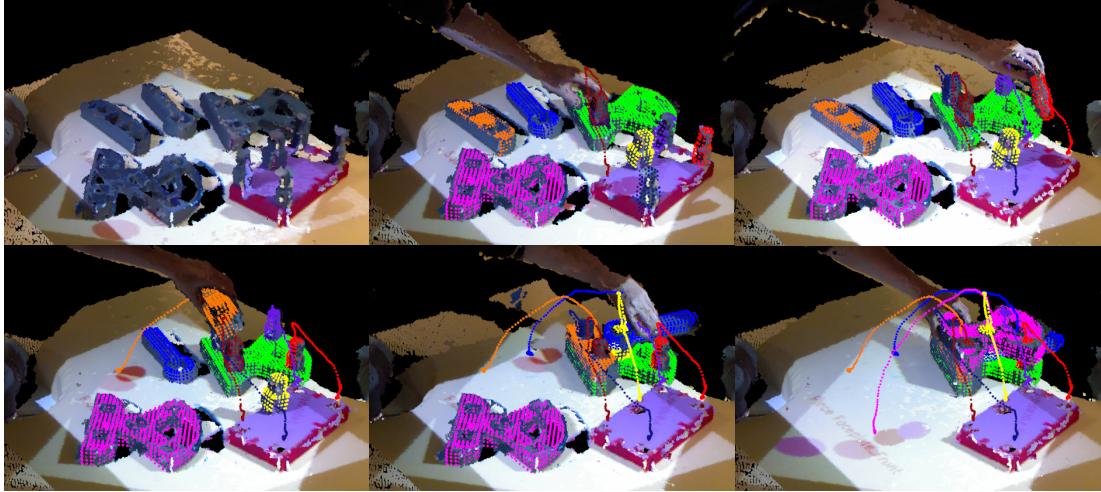


**Figure 4.3.9:** Results on the Tide artificial sequence. Each colored curve represents a certain number of samples, and gives mean RMS error averaged over 25 trial runs for increasing numbers of particles.

### 4.3.2 RESULTS ON REAL SEQUENCES

One application of our tracker is to provide semantic understanding and imitation of assembly tasks. This can be accomplished by tracking all interacting parts of an assembly as a human

demonstrates, and then using the trajectories and poses in order to train a robot to replicate the construction. Additionally, the tracked output can be used as an input for the robot during construction in order to verify that it has successfully completed each step of the task.



**Figure 4.3.10:** Human demonstration of assembly of the Cranfield Scenario. Tracking runs live for all objects at once at sufficient frame rates to track the whole task.

As a demonstration of this, we shall once again use the well established “Cranfield” benchmark set [28], consisting of eight pieces which can be assembled in a number of different orders. In our experiments, models consist of voxelized point clouds derived from high-resolution models of the pieces, and initial poses for tracking are found using a combined object recognition and pose estimation algorithm [22]. Each object is tracked using an independent particle filter, with  $N_{samples}$  set to 50, and  $N_{particles}$  set to 1000.

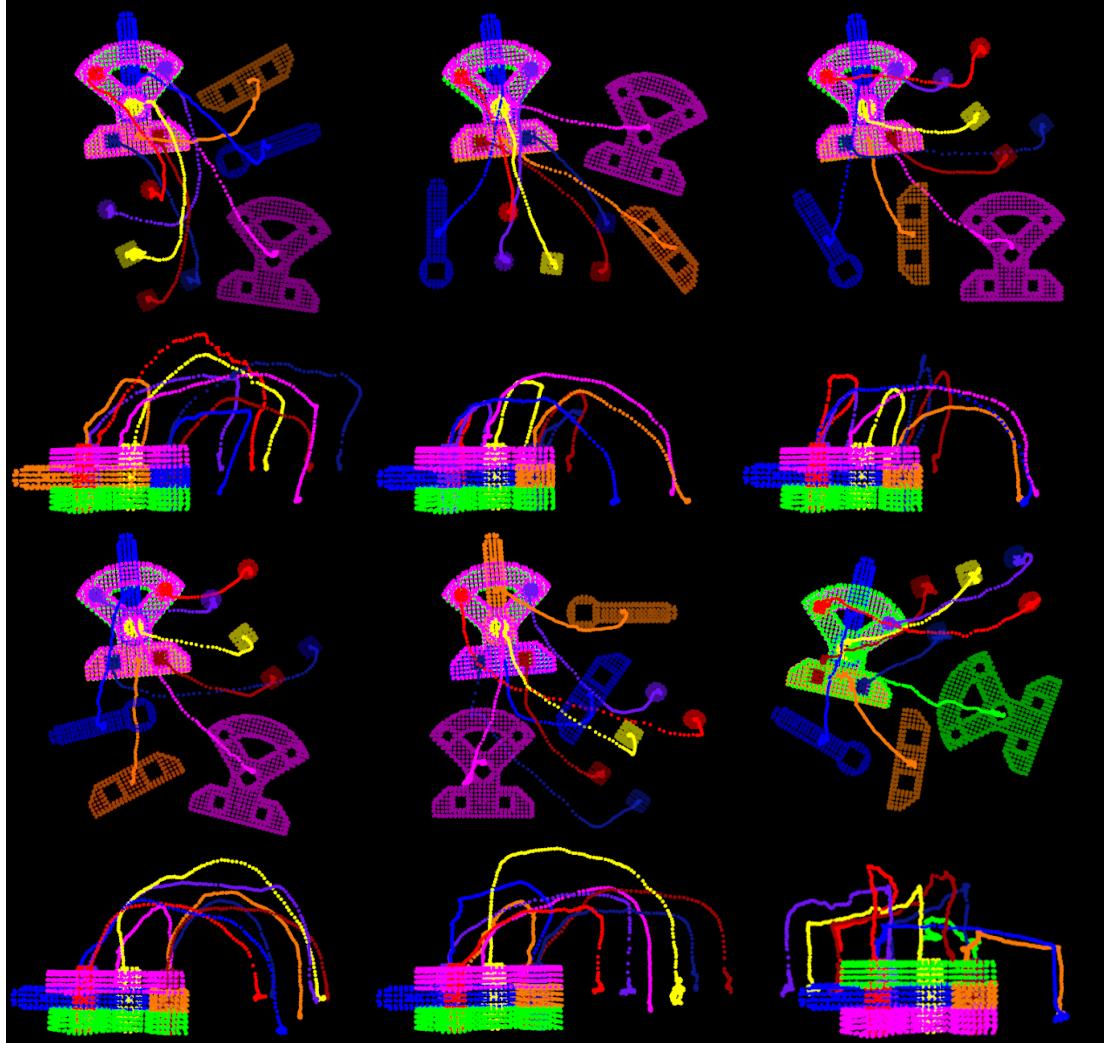
Recordings were made on the MARVIN platform at the University of Southern Denmark, and use 2 Kinect RGB-D cameras<sup>1</sup>. The recordings were performed by different people, where the people were following assembly instructions presented by the planning system of the IntellAct project. A description of the planner is beyond the scope of this work (we refer the reader to [73]), but for our purposes we just need to know that the order of assembly varies from sequence to sequence.

Figure 4.3.10 shows a montage of screenshots captured as a human demonstrates assembly of the benchmark. As can be seen, all pieces are successfully tracked from start to finish, with each tracker outputting smooth trajectories that can be used for training a robot using Dynamic Motion Primitives (DMP) [51]. In Figure 4.3.11 we show tracks from multiple different human demonstrations - one can observe the different strategies that people employ in assembling the benchmark. The tracks in the lower right corner of the Figure are from a robot

---

<sup>1</sup>It is well-known that multiple Kinect sensors sharing a common field of view will cause IR interference, resulting in poor depth reconstructions. A known solution, which the platform incorporates, is the use of vibrating motors mounted on the Kinect sensors [23]. This method has been shown to effectively blur out the noisy contributions of external sensors, while maintaining a high depth reconstruction quality.

reproducing the assembly after being trained on the human demonstrations [73].



**Figure 4.3.11:** Tracking results from six different recordings of the Cranfield Scenario. The tracks in the bottom right corner are from the robot constructing the object, while the other five are from five different human demonstrators. In the overhead views, starting poses are shown (in slightly darker colors) for the objects.

#### 4.4 DISCUSSION

In this Chapter we have presented a novel spatially stratified sampling approach which greatly reduces the computational complexity of 3D Point Cloud correspondence particle filters. We evaluated the tracker using synthetic sequences for which precise ground truth exists, as well as real sequences of a human demonstration application. To demonstrate the effect of stratified sampling on performance, we conducted a sweep over the parameter space of number of particles and samples. This sweep showed the clear effectiveness of the proposed method in matching and even out-performing a GPU implementation.

The approach we have presented here allows us to effectively track rigid objects in 3D voxel space. While it is very efficient at doing this tracking, it remains just that, a tracker, the result of which is an object state. Moreover, it is unable to handle deforming objects, and as of yet we have not shown how to handle objects entering or leaving the scene. That is to say, we have yet to show how we can use this tracking system to produce a full video segmentation.

As such, now that we have now established our ability to track multiple objects in real-time within point cloud data, we can proceed to VOS. In the next Chapter, we will combine the 3D correspondence-based tracker with the supervoxel world-model presented previously to generate a full segmentation of point cloud video that is robust to occlusion and maintains object identities throughout extended sequences. To do this, we shall borrow some ideas from our 2D tracker presented in Chapter 2, and introduce a new global energy function that allows us to assign newly observed supervoxels to tracked targets.

*Encoded in the large, highly evolved sensory and motor portions of the human brain is a billion years of experience about the nature of the world and how to survive in it. The deliberate process we call reasoning is, I believe, the thinnest veneer of human thought, effective only because it is supported by this much older and much more powerful, though usually unconscious, sensorimotor knowledge. We are all prodigious olympians in perceptual and motor areas, so good that we make the difficult look easy.*

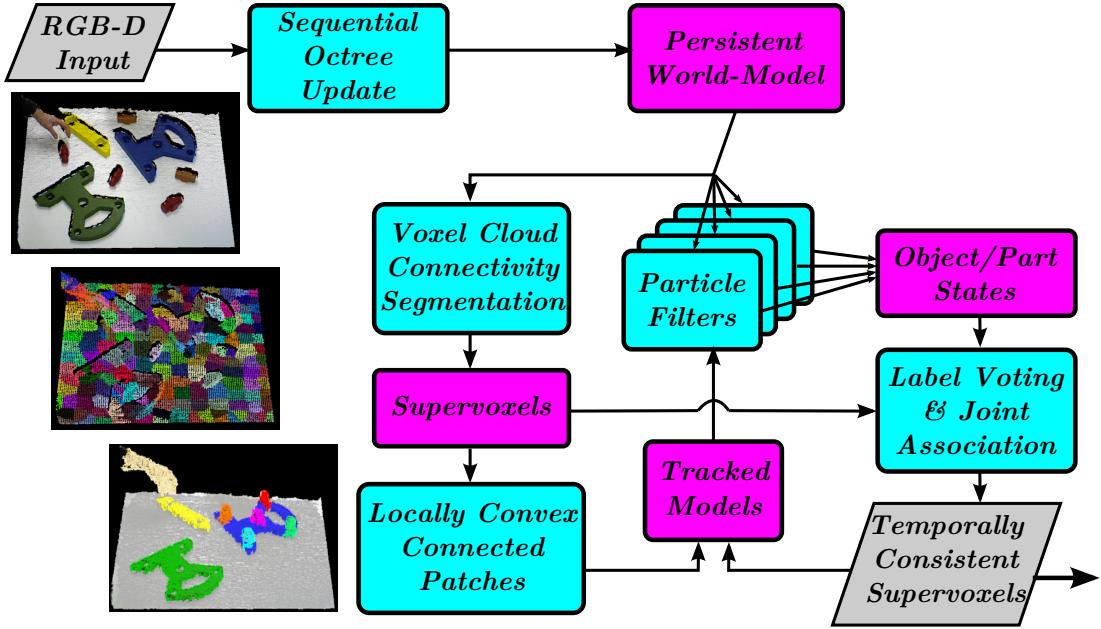
Hans Moravec

# 5

## Tracking Based Point Cloud Video Segmentation

**T**HUS FAR, WE HAVE PRESENTED a 2D particle-filter based VOS method, developed a 3D point-cloud based world-model, and shown how it is possible to efficiently track within this world using correspondence-based particle filters. Our final task is to bridge the gap between the tracked model states (which are the final output from the previous Chapter) and the supervoxels presented in the preceding one. That is to say, we wish to use our tracked results to link supervoxels from frame to frame; essentially, to solve the association problem at the lowest level possible thereby achieving temporally consistent supervoxels.

The motivation for making temporal connections at the supervoxel level (rather than at the significantly easier object level) is to avoid the need to make strict decisions about objects and their boundaries. We wish to avoid these decisions, as what one defines as an “object” is largely dependent on context, as it is really a property of the observer, rather than the observed. By tracking supervoxels we can avoid the problem completely. Instead, we make “fuzzy” associations, where instead of a binary association decision, we instead maintain probabilities that supervoxels “belong” to different tracked entities. An overview of the proposed method is shown in Figure 5.1.1; as can be seen, we use many of the components presented previously, with the main addition being an association step which assigns supervoxels to tracked objects.

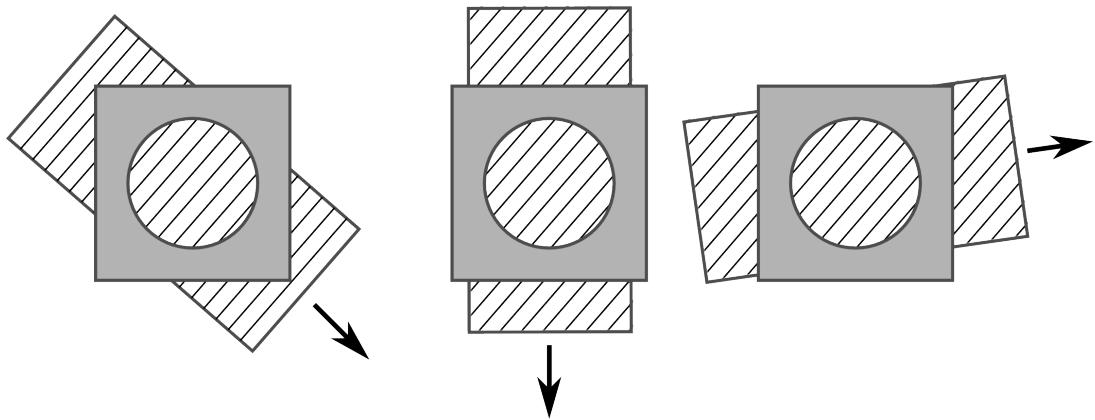


**Figure 5.1.1:** Overview of the algorithm for extracting temporally consistent supervoxels. The persistent world model, VCCS, LCCP, and particle filters function as presented in previous Chapters. The key addition is the label voting and joint association scheme which uses tracked states to associate supervoxels from frame to frame. The output of this is then fed back to the trackers to update their models.

## 5.1 TRACKED MODEL REPRESENTATION

The first issue that must be addressed is the “level” at which targets should be tracked - the object or the supervoxel level (we shall not consider tracking at the voxel level since it is computationally infeasible with current hardware). As our goal is to associate supervoxels across time, we would like to track supervoxels directly. Unfortunately, this is generally not feasible due to the “aperture problem” seen in neural visual fields [56]. The aperture problem deals with the fact that local motion can only be estimated perpendicular to a contour that extends beyond its field of view [78]. In other words, determining direction of motion in a local region (without considering global features) is generally not possible - as illustrated in Figure 5.1.2. This means that in order to estimate motion of supervoxels, we must extend the field of view considered significantly beyond the size of the supervoxel itself; in fact, our aperture must contain the borders of the moving object in question, otherwise pairwise association of supervoxels is generally indeterminate.

Thud we must track higher level groupings - groupings that extend at least to a contour which provides a reference boundary for disambiguating motion. A natural way of doing this is to use the LCCP segmentation presented in Chapter 3, as it will expand regions up to concave boundaries. Using concave connections as references is surprisingly powerful, as it generally



**Figure 5.1.2:** The “Aperture Problem” - Three patterns moving in three different directions all produce the same perceived stimulus when viewed through a small aperture. In order to correctly resolve direction of motion a wider field of view is needed. Adapted from [45].

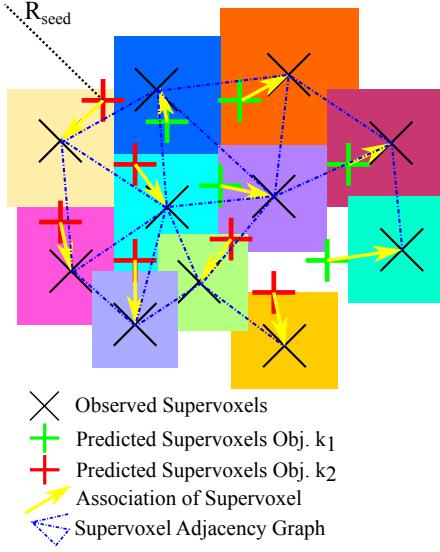
will differentiate objects, as well as parts of objects which can move independently (consider the case of joints in the human body). As such, we adopt a simple scheme for grouping supervoxels into tracked models; we perform LCCP segmentation on the first frame, and assign each observed segment to an independent tracker.

## 5.2 BANK OF PARALLEL PARTICLE FILTERS

Tracking of the segmented objects or parts is accomplished using a bank of the correspondence-based particle filters from the previous Chapter. We select a measurement model based on the voxels and supervoxels produced using the persistent world-model scheme discussed in Chapter 3. The model measures distance in a feature space of spatial distance, normals, and color. Weights of predicted states  $\mathbf{x}_t^j = [d_x, d_y, d_z, \gamma, \beta, a]$  are measured as in the previous Chapter by associating transformed model voxels to the observed voxels nearest in space. Particles are then weighted by measuring total distance in feature space, just as in Equation 4.10.

## 5.3 ASSOCIATION BY JOINT LABEL OPTIMIZATION

Since our goal is extract full segmentations, rather than just object states, we must actively associate observed supervoxels (and therefore, voxels) with tracked results. This can be considered as an additional step beyond the tracking of the previous Chapter, which used observations to test how well a particle prediction matched reality. We can begin by considering the trivial case of tracking a single object. In this case we can extract a segmentation by establishing a small search radius  $R_{assoc}$  around each tracked model result point (giving us an “association volume” around each point). If we simply assign each observed voxel falling within this vol-



**Figure 5.3.1:** Association of observed supervoxels with predicted model supervoxels using smoothing term which considers neighbor labels.

ume to the tracked model, we can easily achieve a rough segmentation. Furthermore, we can determine which supervoxels belong to the object (and which belong to the unlabeled set) using a majority-voting scheme. This gives us a “foreground segmentation” method, which segments out the tracked object as foreground.

Extending this to multiple objects, observe that the difficulty now lies in resolving associations of voxels which lie within the association volume of two or more tracked results. Fortunately, most voxels will only fall within the volume of a single object’s tracked model, leaving us with only having to resolve associations along interacting object boundaries. As such, we begin by composing a list of all supervoxels which are under competition, that is, have voxels falling within the association volumes of more than one label. We can then count the number of voxels associated to each object, and normalize to give us an a-priori categorical distribution  $P(L(p) = k|V)$  which maps labeling of supervoxels  $p \in P$  to objects  $k \in K$  given voxel associations  $V$ .

Now that we have priors for object labelings, we adopt a Monte Carlo approach, similar to [43], to sample from the set of possible label associations and determine a global association which best aligns tracked object predictions to observed supervoxels. To generate realizations (sets of label assignments), we use a weighted sampling strategy which considers the priors, as well as a distance term. This gives us a likelihood of assigning object label  $k$  to supervoxel  $p$  given distance from the object centroid  $C_k$  and the voxel associations:

$$\mathcal{L}(L(p) = k|C_k, V). \quad (5.1)$$

To compute a score for each realization, we use the global energy function given in (5.2).

Each global label association  $\mathcal{A}$  consists of a set of associations  $\{a_1 \dots a_n\}$  which assign object labels  $k$  to the set of observed supervoxels  $\{p_1 \dots p_n\}$ . The first summation term,  $\sum_p \|p_k^* - p\|$ , measures error in feature space between the observed supervoxel and the supervoxel of the stratum in its associated object  $p_k^*$ .

$$E_{\mathcal{A}} = \left( \sum_p \|p_k^* - p\| + \lambda \sum_{(p,p') \in \mathcal{N}} \delta(L(p) \neq L(p')) \right) \prod_{a \in \mathcal{A}} \Delta_k \quad (5.2)$$

The second summation is a smoothing term which considers the adjacency graph of observed supervoxels. For every observed supervoxel, we compare its assigned labeling  $L(p)$  to the label of all supervoxels  $p'$  which lie within its adjacency neighborhood  $\mathcal{N}$ . We adopt the Potts model as in [19], where  $\delta(\cdot)$  is 1 if the specified condition holds, and 0 otherwise, and  $\lambda$  is a weighting coefficient which controls the importance given to spatial smoothness of labels.

Finally, the multiplicative term  $\prod_{a \in \mathcal{A}} \Delta_k$  controls for the expansion or contraction of objects through the number of observed voxels associated with them.  $\Delta_k$  penalizes for changes in volume by increasing the energy for deviations from unity in the ratio of observed voxels assigned to an object  $\hat{N}_k$  with the number in the object model itself  $N_k$ , that is

$$\Delta_k = \begin{cases} \hat{N}_k/N_k & \text{if } \hat{N}_k \geq N_k \\ 2 - \hat{N}_k/N_k & \text{if } \hat{N}_k < N_k . \end{cases} \quad (5.3)$$

Once we have arrived at a stable minimum energy score, we extract the resulting association of observed supervoxels to predicted results, and use them to update the tracked models.

#### 5.4 ALIGNMENT AND UPDATE OF MODELS

The joint energy minimization results in a global association  $\mathcal{A}$  which assigns observed supervoxels to tracked objects. In order to use this to update the object models, we must align it to the internal representation stored by the particle filter. We begin at the inverse of the predicted state, and then use an iterative closest point [25] procedure to refine the transform such that the set of observed supervoxels best aligns with the model. We then update the model with the new observations by inserting the new supervoxels into the model.

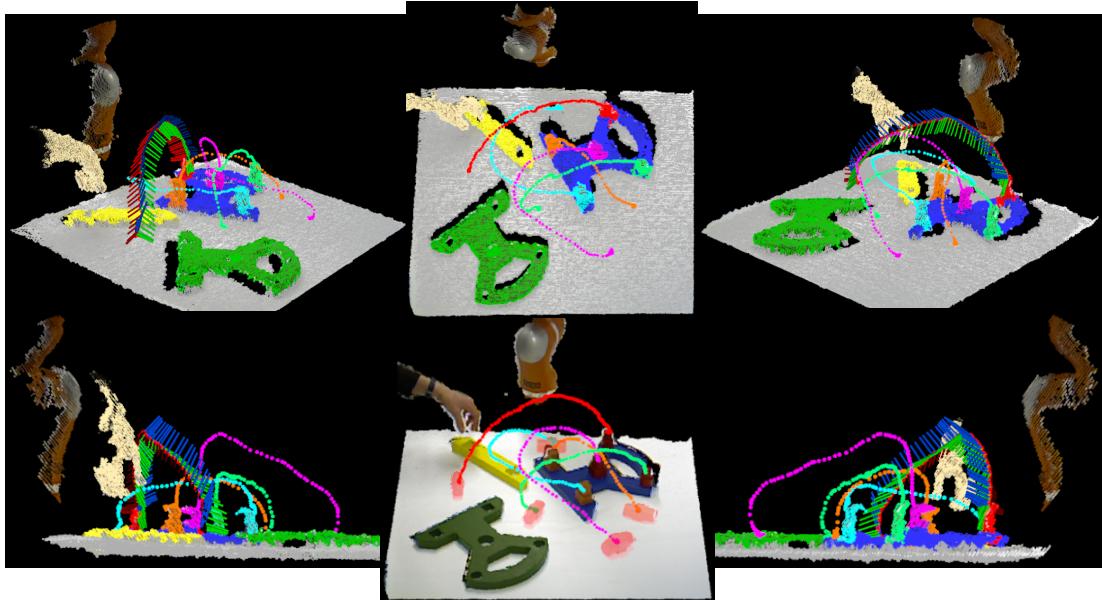
As a final step, we use the refined transform to update the states of the particles. To do this, we shift each particle  $x_i$  towards the refined state  $\hat{x}$ , weighting the importance given to the refined state by a constant factor  $\varepsilon$

$$x'_{i \in L} = (1 - \varepsilon)x_i + \varepsilon\hat{x} . \quad (5.4)$$

For this work, we found that an  $\varepsilon$  of 0.5 effectively removes noise (jitter) introduced by the replacement of the tracked model. Additionally, we correct the internal motion model ( $\{\nu_x, \nu_y, \nu_z\}$ ) of the particle filters to correspond to the new updated state.

## 5.5 EXPERIMENTAL RESULTS

As a demonstration of the method, in this Section we provide results from two successful applications. Both applications use the Cranfield scenario [27] presented in the previous Chapter. Figure 5.5.1 shows the results of tracking and segmentation on one assembly of the benchmark. It can be seen that the algorithm is able to successfully extract full segmentation of the video, as seen by the tracks and the segmented pieces.

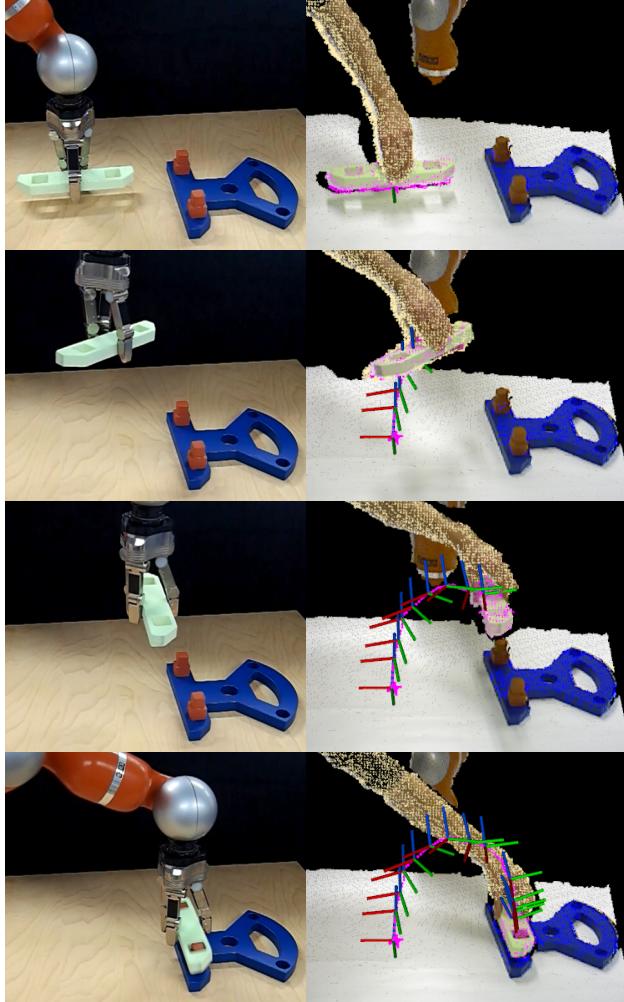


**Figure 5.5.1:** Result of tracking and segmentation on Cranfield scenario from different views. Here the tracks are shown as dots of the color of the tracked label for each timestep. Initial locations of the pegs are shown in the middle bottom frame as semi-transparent masks. Calculated orientation is shown for the red peg with a set of axes every second time-step; these axes show pose in a frame relative to the start.

### 5.5.1 IMITATION OF TRAJECTORIES FOR ROBOT MANIPULATION

The standard way of teaching robots to perform human-like actions is imitation learning, also called programming by demonstration [13, 16]. There are several ways to demonstrate movements: 1) recording movements in joint-space (joint angles) or target-space (Cartesian space) by ways of a motion capture device (requires putting markers on human body), 2) using kinesthetic guidance (guiding a robot's movements by a human hand), or 3) via teleoperation (controlling a robot via joystick). The only way to obtain motion trajectories from human observation in a "non-invasive" procedure is by using stereo vision [38], however, usually it is model based. The tracking algorithm we have presented here can be used as an alternative method to obtain motion trajectories (in Cartesian space) in a model-free way.

To demonstrate this, we applied our tracking algorithm to obtain human motion trajec-

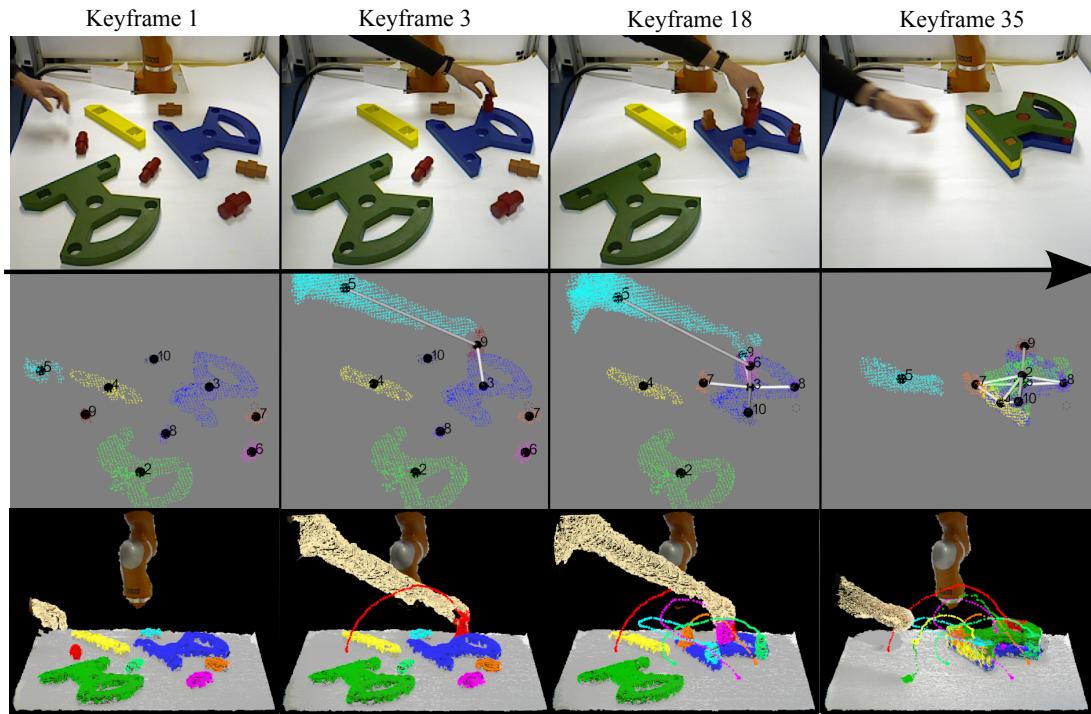


**Figure 5.5.2:** KUKA LWR arm imitating trajectory and pose learned from tracked human demonstration.

ries in Cartesian space including orientation of manipulated object (in total six DoFs). We tested it using a recording of the Cranfield scenario where, first, we let a human demonstrate the action and then reproduced it using a KUKA Light Weight Robot (LWR) arm [50]. Specifically, here we imitate a human putting the separator block on the pegs. To generate trajectories for the robot from human demonstrations, we used a modified version of Dynamic Movement Primitives [41, 42] (DMP) and learning method as described in [51]. We used Cartesian impedance control and, thus, generated six DMPs (three for motion of the end-effector in Cartesian space and three for orientation of the hand) based on trajectories obtained from the tracking algorithm. Here we used 100 equally spaced kernels with width  $\sigma = 0.05$  for each dimension (for more details please refer to [51]). As demonstrated in Figure 5.5.2 and the supplementary video, trajectories obtained by the proposed tracking algorithm are sufficiently accurate to allow reproduction of the human motion. We should emphasize that the key advantage here over the tracking from the previous Chapter is that we can track and segment out the human arm - something not possible with a rigid-model approach.

### 5.5.2 SEMANTIC SUMMARIES OF ACTIONS

A fundamental task for intelligent autonomous robots is the problem of encoding long chain manipulations in a generic way, for use in tasks such as learning and recognition. As a demonstration of the usefulness of the proposed tracking framework, we use a recently introduced novel Semantic Event Chain (SEC) approach [9] which converts each segmented scene to a graph: nodes represent segment (i.e. object) centers and edges indicate whether two objects touch each other or not. By using an exact graph matching technique the SEC framework discretizes the entire graph sequence into decisive main graphs. A new main graph is identified whenever a new node or edge is formed or an existing edge or node is deleted. Thus, each main graph represents a “key frame” in the manipulation sequence. Figure 5.5.3 shows a few detected sample key frames from the long Cranfield action. While the complete action has in total 1453 frames, the SEC representation reduces it to just 12 key frames, each of which represents a topological change in the scene.



**Figure 5.5.3:** A few example key frames extracted from the long Cranfield action. Numbered nodes represent interacting objects, while edges show touching relations between objects. Each keyframe represents a topological change in the scene - here we show 4 of the 12 keyframes.

## 5.6 DISCUSSION

In this Chapter we presented a method which extracts full scene segmentation from the model tracker we presented in the previous Chapter. The method uses a global energy functional to enforce smoothness and temporal continuity on the segmentations. Additionally, we feed

segmentation results back into the tracked models, updating them so that they can deform and/or accumulate different view points. This allows us to extract a full video segmentation from our tracked states.

There are many advantages to the method presented in this Chapter over the model tracker - it can track deforming objects, there is no need to know object models a-priori, tracked objects can be easily re-initialized using LCCP segmentation, there is no need to compute initial poses, and finally, we can extract action semantics without needing to know what objects are present. As this approach is completely free of trained models or strict object boundaries, it opens up many avenues of future research. In particular, of especial interest is that it allows bootstrapping of learning - we can attempt to build systems which learn to understand scenes purely from observation, without any human input or teaching.



*Art thou not, fatal vision, sensible  
To feeling as to sight? Or art thou but  
A dagger of the mind, a false creation  
Proceeding from the heat-oppressed brain?*

William Shakespeare

# 6

## Conclusions

**T**HROUGHOUT THIS WORK, we have had one goal in mind; to develop video segmentation which has the continuity of tracking methods. The primary reason for doing this was to ensure the temporal consistency of segments through extended video clips of, in particular, indoor manipulation tasks. Difficulties presented in such videos include partial and full occlusions, sudden and fast displacements of objects, different objects with similar or identical appearance, and objects which cannot be segmented based on color. Additionally, proper understanding of manipulations require a high degree of precision, particularly when it comes to the relative pose of parts when they are interacting (e.g. putting a bolt in a hole).

Our intended application for the segmentation of such videos was the general bootstrapping of assembly task understanding. If one can correctly track objects and their parts through a task without a-priori knowledge, it should be possible to use this to learn without the need of an external oracle. Not only this, but if one is able to correctly track full 6 DoF pose throughout a human-demonstrated assembly task, it is possible to learn trajectories that a robot can use to imitate effective motion paths.

### 6.1 SUMMARY OF CONTRIBUTIONS

We began in chapter 2 by presenting a 2D VOS method that made use of our proposed methodology; to use tracking as the basis for video segmentation. In particular, we showed how parti-

cle filters, a class of Sequential Bayesian Monte Carlo methods, can be used to predict what the next frame's segmentation should look like. These proposed segment masks were then combined using a weighted sampling strategy and then refined to fit observed image data using a relaxation process.

Next, in chapter 3 we introduced RGB-D sensors, how they can be used to produce 3D point cloud data, and how this data can be organized using an octree structure. We then presented a specialized type of octree - the adjacency octree - which we developed to allow quick and efficient traversal between neighboring voxels within the tree. We subsequently showed how the adjacency octree can be used to efficiently further sub-divide voxel data into localized patches, called supervoxels, using our Voxel Cloud Connectivity Segmentation (VCCS). The utility of supervoxels was then demonstrated by showing their effectiveness in segmenting static scenes using a local convexity criterion (LCCP). The effectiveness of VCCS and LCCP were then demonstrated by showing their favorable results on a large state-of-the-art benchmark as compared to other state-of-the-art methods. Finally, we conclude the chapter by presenting how an adjacency octree containing supervoxels can be sequentially updated with new frames of data without deleting potentially occluded voxels.

We then extended the particle filter framework in chapter 4 to 3D point clouds by formalizing the notion of a voxel-based measurement and dynamic model. While this straightforward implementation worked, we showed how it could achieve much faster (real-time on standard hardware) run times and accuracy by stratified sampling of correspondences. This improvement was then quantified using a benchmark of artificial sequences, and then shown to outperform even a GPU based method. As a final demonstration of the effectiveness of the tracker, we presented results on recordings of humans constructing the Cranfield benchmark. This showed how the tracker can be used to distill semantic understanding from a video sequence.

Finally, in chapter 5 we tackled the problem of extracting full segmentations from tracked results. To do this, we first showed how the 3D particle filter presented previously could be extended to work on supervoxels. Then we showed how the adjacency graph of supervoxels could be used along with a global energy minimization to resolve interactions between trackers and produce a full segmentation consistent with the tracked poses. As a final demonstration of the presented methodology, we give results on several recordings of manipulation actions.

Another important contribution was the development of the open-source Oculus vision system discussed in Appendix A. This system served as the platform on which much research has been done over the past several years and was a key tool in publications by several other researchers. Finally, we would like to note that most of the algorithms discussed in this work have been released as open-source to the vision community as part of the Point Cloud Library<sup>1</sup>. We consider both of these important contributions, as the open sharing of code is vital to the advancement of the discipline of Computer Science.

---

<sup>1</sup><http://www.pointclouds.org/>

## 6.2 SHORTCOMINGS OF VOS BENCHMARKS

Evaluation of segmentation algorithms is a notoriously vexing problem due to the inherent ambiguity of what constitutes a “correct” segmentation. As such, in our work we have determined that we should avoid making concrete decisions on segmenting objects in single works, and instead chose to limit ourself to the lower, supervoxel level. While this is not an entirely satisfactory solution, we felt that there is simply not enough information in a single frame to extract meaningful segmentations accurately. Indeed, one cannot really tell the granularity with which a scene should be segmented into distinct objects until they see some action.

Benchmarking of 3D point cloud segmentation is a young topic - in fact, the first extensive benchmark, the *NYU Indoor Dataset*, was not published until 2012 [79]. As such, it has many complications (that did not exist in 2D) which have yet to be resolved, such as that it is difficult, if not impossible, to make a 2D ground truth annotation correctly line up with the 3D point cloud representation. Even more to the point, there are still no 3D VOS benchmarks. In fact, even though the field is decades old, one must look to 2013 to find a 2D VOS benchmark [34]. There are many reasons for this lack of a proper benchmark, but the primary one is that is simply extremely time consuming to annotate ground-truth for even very short video sequences. Furthermore, labeling a single ground truth is even more difficult for video than single images, for instance, what happens when one takes a cap off of a bottle; should it be given a new label? If so, should it have had a separate label the entire time, or only once it is separated? What happens when objects become occluded and then reappear; should they be given new identities or maintain their old ones? If they keep their old ones, how long should we allow an object to be occluded for before we “forget” it?

Due to all of these concerns, we have made the decision to only show qualitative results of segmentation per-pixel accuracy, while still quantifying the tracker performance. This allowed us to prove that our tracking scheme was both faster and more accurate than existing methods, without needing to haggle over inscrutable questions such as “what constitutes an object?”, and “does this pixel belong to object a or object b?”.

## 6.3 LIMITATIONS AND DIRECTION OF FUTURE WORK

The main limitation of the 2D tracking framework presented in chapter 2 is that it can only “guess” at correct behavior when an object is occluded. Indeed, this is a general problem of trying to infer behavior in a 3D world from observations in a 2D projected plane. It is because of these ambiguities and an inability to resolve them in a comprehensive and satisfying manner that we proceeded to tracking in 3D using RGB+D observations.

While our persistent voxel world model is very effective at maintaining the existence of stationary objects through occlusions, it does not handle objects which move while they are occluded. Solving this problem at the low-level of voxels is an unresolved problem, the outlook

of which is fairly bleak. Our attempts at solving the problem lead us to believe that higher level object knowledge is necessary to account for occluded motion. With this in mind we are investigating a way of associating occluded objects with their occluder so that they move with them. Another limitation which we are currently addressing is that our persistent voxel world model does not account for camera motion. There is some existing work on real-time camera pose estimation, and we are hoping to incorporate such a method into our system in the near future.

As with any VOS method, our final result has a few important limitations. One of these is the need to set a rate at which to allow models to change. In many cases, objects are rigid, and do not change; allowing them to change only adds instability to the segmented output. Conversely, some objects are deformable, or not entirely visible in the scene, and will need to change to be correctly segmented.

Future work should focus on extending how patch-based VOS is extracted from the trackers. While we have shown that using tracking as the basis for video segmentation is a viable concept, there remains some work to be done in terms of reliability of the tracked patches. In particular, one promising avenue to pursue is using a global Conditional Random Field to help give better solutions to the joint association of supervoxels to trackers.

Finally, while it was beyond the scope of this work, the clear next step is to take the tracked patches and use them to bootstrap learning from visual data. That is, now that we have established that it is possible, in principle, to track local patches without the need for a-priori object knowledge, we can begin to learn objects from unconstrained video. Possible places to begin such research would be to attach relevant metadata to videos to use for semi-supervised learning, or to combine different sensing modalities. For instance, one could combine speech commands and a video showing an execution of the action into a single data structure. One could then use tracked patches to discover objects related to the action, and by observing similarities over multiple instances of similar actions, correlate video observations with speech - allowing discovery of what objects correspond to what vocal commands, or how ordering of spoken commands corresponds to observed actions.

We propose that this is how the field shall advance, how true artificial cognitive agents can be created; by allowing unsupervised learning on video data. With this in mind, we have attempted to reduce video data from an incomprehensible mess of discrete, disconnected pixels into a graph of temporally continuous patches. Importantly, we have avoided making “object” decisions - something we believe must be determined through observations, and cannot be encoded in rules on a set of features. By tracking patches, we have created a representation of video which both structures and reduces the raw sensory input enough that it can be used for learning while remaining flexible to unknown objects. Our methods, and the representation they produce, give local structure as well as temporal continuity at the lowest-level of visual perception - establishing the structure on which one can build high-level understanding.

# Bibliography

- [1] V. Ablavsky, A. Thangali, and S. Sclaroff. Layered graphical models for tracking partially-occluded objects. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [2] A. Abramov, K. Pauwels, J. Papon, F. Worgotter, and B. Dellen. Real-time segmentation of stereo videos on a portable system with a mobile gpu. *Circuits and Systems for Video Technology, IEEE Transactions on*, Sept 2012.
- [3] A. Abramov, K. Pauwels, J. Papon, F. Worgotter, and B. Dellen. Depth-supported real-time video segmentation with the kinect. In *Applications of Computer Vision (WACV), 2012 IEEE Workshop on*, Jan 2012.
- [4] Alexey Abramov, Eren Erdal Aksoy, Johannes Dörr, Florentin Wörgötter, Karl Pauwels, and Babette Dellen. 3d semantic representation of actions from efficient stereo-image-sequence segmentation on gpus. In *International Symposium 3D Data Processing, Visualization and Transmission*, 2010.
- [5] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Trans. Pattern Anal. Machine Intell.*, 34(11):2274–2282, nov. 2012. ISSN 0162-8828. doi: 10.1109/TPAMI.2012.120.
- [6] A. Adam, E. Rivlin, and I. Shimshoni. Robust fragments-based tracking using the integral histogram. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006.
- [7] N. Ahuja and S. Todorovic. Connected segmentation tree; a joint representation of region layout and hierarchy. In *CVPR*, pages 1–8, 2008.
- [8] E. E. Aksoy, A. Abramov, F. Wörgötter, and B. Dellen. Categorizing object-action relations from semantic scene graphs. In *ICRA*, pages 398–405, may 2010.
- [9] E. E. Aksoy, A. Abramov, J. Dörr, K. Ning, B. Dellen, and F. Wörgötter. Learning the semantics of object-action relations by observation. *The International Journal of Robotics Research*, 30(10):1229–1249, 2011.
- [10] E. E. Aksoy, A. Abramov, J. Dörr, K. Ning, B. Dellen, and F. Wörgötter. Learning the semantics of object-action relations by observation. *The International Journal of Robotics Research*, 30(10):1229–1249, 2011.
- [11] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE Trans. PAMI*, 33(5):898–916, 2011.

- [12] P. Arbelaez, B. Hariharan, Chunhui Gu, S. Gupta, L. Bourdev, and J. Malik. Semantic segmentation using regions and parts. In *CVPR*, pages 3378–3385, 2012.
- [13] B. Argall, S. Chernova, M. M. Veloso, and B. Browning. A survey of robot learning from demonstration. *Robots and Auton. Sys.*, 57(5):469–483, 2009.
- [14] B. Babenko, Ming-Hsuan Yang, and S. Belongie. Visual tracking with online multiple instance learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [15] A. Barla, F. Odone, and A. Verri. Histogram intersection kernel for image classification. In *Image Processing, 2003. ICIP 2003. Proceedings. 2003 International Conference on*, pages III – 513–16 vol.2, sept. 2003. doi: 10.1109/ICIP.2003.1247294.
- [16] A. Billard, S. Calinon, R. Dillmann, and S. Schaal. *Survey: Robot Programming by Demonstration*. MIT Press, 2008.
- [17] M. Blatt, S. Wiseman, and E. Domany. Superparamagnetic clustering of data. *Physical Review Letters*, 76(18):3251–3254, 1996.
- [18] L. Bourdev and J. Malik. Poselets: Body part detectors trained using 3D human pose annotations. In *ICCV*, pages 1365–1372, 2009.
- [19] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Machine Intell.*, 23(11):1222–1239, Nov 2001. ISSN 0162-8828. doi: 10.1109/34.969114.
- [20] M.D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier, and L. Van Gool. Online multiperson tracking-by-detection from a single, uncalibrated camera. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(9):1820–1833, Sept 2011.
- [21] W. Brendel and S. Todorovic. Video object segmentation by tracking regions. In *IEEE International Conference on Computer Vision (ICCV)*, 2009.
- [22] Anders Glent Buch, Yang Yang, Norbert Krüger, and Henrik Gordon Petersen. In search of inliers: 3d correspondence by local and global voting. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, June 2014.
- [23] D Alex Butler, Shahram Izadi, Otmar Hilliges, David Molyneaux, Steve Hodges, and David Kim. Shake’n’sense: reducing interference for overlapping structured light depth cameras. In *Proceedings of the 2012 ACM annual conference on Human Factors in Computing Systems*, pages 1933–1936. ACM, 2012.
- [24] Yizheng Cai, Nando Freitas, and JamesJ. Little. Robust visual tracking for multiple targets. In Aleš Leonardis, Horst Bischof, and Axel Pinz, editors, *Computer Vision – ECCV 2006*, volume 3954 of *Lecture Notes in Computer Science*, pages 107–118. Springer Berlin Heidelberg, 2006. doi: 10.1007/11744085\_9.
- [25] Dmitry Chetverikov, Dmitry Stepanov, and Pavel Krsek. Robust euclidean alignment of 3d point sets: the trimmed iterative closest point algorithm. *Image and Vision Computing*, 23(3):299 – 309, 2005. ISSN 0262-8856. doi: 10.1016/j.imavis.2004.05.007.

- [26] Changhyun Choi and H.I. Christensen. Rgb-d object tracking: A particle filter approach on gpu. In *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, Nov 2013. doi: 10.1109/IROS.2013.6696485.
- [27] K Collins, AJ Palmer, and K Rathmill. The development of a european benchmark for the comparison of assembly robot programming systems. In *Proceedings of the 1st Robotics Europe Conference, Brussels*, pages 27–28, 1984.
- [28] K. Collins, A. J. Palmer, and K. Rathmill. The development of a European benchmark for the comparison of assembly robot programming systems. In *Robot technology and applications (Robotics Europe Conference)*, pages 187–199, 1985.
- [29] D. Comaniciu, V. Ramesh, and P. Meer. Real-time tracking of non-rigid objects using mean shift. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2000.
- [30] Arnaud Doucet, Nando De Freitas, and Neil Gordon, editors. *Sequential Monte Carlo methods in practice*. 2001.
- [31] P.F. Felzenszwalb, R.B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE Trans. PAMI*, 32(9):1627–1645, 2010.
- [32] T.E. Fortmann, Y. Bar-Shalom, and M. Scheffe. Multi-target tracking using joint probabilistic data association. In *Decision and Control including the Symposium on Adaptive Processes, 1980 19th IEEE Conference on*, volume 19, pages 807–812, Dec 1980.
- [33] Thomas E. Fortmann, Y. Bar-Shalom, and M. Scheffe. Sonar tracking of multiple targets using joint probabilistic data association. *Oceanic Engineering, IEEE Journal of*, 8(3):173–184, Jul 1983.
- [34] Fabio Galasso, Naveen Shankar Nagaraja, Tatiana Jiménez Cárdenas, Thomas Brox, and Bernt Schiele. A unified video segmentation benchmark: Annotation, metrics and analysis. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, 2013.
- [35] T. Gautama and M. Van Hulle. A phase-based approach to the estimation of the optical flow field using spatial filtering. *IEEE Transactions on Neural Networks*, pages 1127–1136, 2002.
- [36] N.J. Gordon, D.J. Salmond, and A F M Smith. Novel approach to nonlinear/non-gaussian bayesian state estimation. *Radar and Signal Processing, IEE Proceedings F*, 140(2):107–113, Apr 1993. ISSN 0956-375X.
- [37] Saurabh Gupta, Pablo Arbelaez, and Jitendra Malik. Perceptual organization and recognition of indoor scenes from RGB-D images. In *CVPR*, pages 564–571, 2013.
- [38] F. Hecht, P. Azad, T. Asfour, and R. Dillmann. Markerless human motion tracking with a flexible model and appearance learning. In *Robotics and Automation (ICRA), 2009 IEEE International Conference on*, 2009.
- [39] V. Hedau, H. Arora, and N. Ahuja. Matching images under unstable segmentations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, june 2008.

- [40] C. Hue, J.-P. Le Cadre, and P. Perez. Tracking multiple objects with particle filtering. *IEEE Transactions on Aerospace and Electronic Systems*, 38(3):791 – 812, jul 2002.
- [41] A. Ijspeert, J. Nakanishi, P. Pastor, H. Hoffmann, and S. Schaal. Dynamical movement primitives: learning attractor models from motor behaviors. *Neural Comput.*, (25):328–373, 2013.
- [42] J. A. Ijspeert, J. Nakanishi, and S. Schaal. Movement imitation with nonlinear dynamical systems in humanoid robots. In *Robotics and Automation (ICRA), 2002 IEEE International Conference on*, pages 1398–1403, 2002.
- [43] Hossam Isack and Yuri Boykov. Energy-based geometric multi-model fitting. *International Journal of Computer Vision*, 97:123–147, 2012. ISSN 0920-5691. doi: 10.1007/s11263-011-0474-7.
- [44] Michael Isard and Andrew Blake. Condensation—conditional density propagation for visual tracking. *International Journal of Computer Vision*, 29(1):5–28, 1998.
- [45] Eric R Kandel, James H Schwartz, Thomas M Jessell, et al. *Principles of neural science*, volume 4. McGraw-Hill New York, 2000.
- [46] Zia Khan, T. Balch, and F. Dellaert. Mcmc-based particle filtering for tracking a variable number of interacting targets. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 27(11):1805–1819, Nov 2005.
- [47] Kourosh Khoshelham and Sander Oude Elberink. Accuracy and resolution of kinect depth data for indoor mapping applications. *Sensors*, 12(2):1437–1454, 2012. ISSN 1424-8220. doi: 10.3390/s120201437.
- [48] Seongyong Koo, Dongheui Lee, and Dong-Soo Kwon. Multiple object tracking using an rgb-d camera by hierarchical spatiotemporal data association. In *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, Nov 2013. doi: 10.1109/IROS.2013.6696489.
- [49] J.J. Kuffner. Effective sampling and distance metrics for 3d rigid body path planning. In *Robotics and Automation (ICRA) 2004. IEEE International Conference on*, April 2004. doi: 10.1109/ROBOT.2004.1308895.
- [50] Kuka Robot Systems. URL <http://www.kuka-robotics.com>.
- [51] T. Kulvicius, K. J. Ning, M. Tamasiunaite, and F. Wörgötter. Joining movement sequences: Modified dynamic movement primitives for robotics applications exemplified on handwriting. *IEEE Trans. Robot.*, 28(1):145–157, 2012.
- [52] O. Lanz. Approximate bayesian multibody tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(9):1436–1449, Sept 2006.
- [53] Anat Levin, Dani Lischinski, and Yair Weiss. Colorization using optimization. *ACM Trans. Graph.*, 23(3):689–694, 2004. ISSN 0730-0301. doi: 10.1145/1015706.1015780.
- [54] A. Levinshtein, A. Stere, K.N. Kutulakos, D.J. Fleet, S.J. Dickinson, and K. Siddiqi. Turbopixels: Fast superpixels using geometric flows. *IEEE Trans. Pattern Anal. Machine Intell.*, 31(12):2290 – 2297, dec. 2009. ISSN 0162-8828. doi: 10.1109/TPAMI.2009.96.

- [55] Siying Liu, Guo Dong, Chye Hwang Yan, and Sim Heng Ong. Video segmentation: Propagation, validation and aggregation of a preceding graph. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [56] D. Marr and S. Ullman. Directional selectivity and its use in early visual processing. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 211(1183):151–180, 1981. doi: 10.1098/rspb.1981.0001.
- [57] Dennis Mitzel and Bastian Leibe. Taking mobile multi-object tracking to the next level: People, unknown objects, and carried items. In *Computer Vision – ECCV 2012*, volume 7576 of *Lecture Notes in Computer Science*, pages 566–579. Springer Berlin Heidelberg, 2012. ISBN 978-3-642-33714-7.
- [58] N. Papadakis and A. Bugeau. Tracking with occlusions via graph cuts. *IEEE Trans. Pattern Anal. Machine Intell.*, 33(1):144–157, Jan. 2011.
- [59] Jeremie Papon and Florentin Wörgötter. Spatially stratified correspondence sampling for real-time point cloud tracking. In *Applications of Computer Vision (WACV), 2015 IEEE Conference on (SUBMITTED)*, jan. 2015.
- [60] Jeremie Papon, Alexey Abramov, Eren Aksoy, and Florentin Wörgötter. A modular system architecture for online parallel vision pipelines. In *Applications of Computer Vision (WACV), 2012 IEEE Workshop on*, jan. 2012.
- [61] Jeremie Papon, Alexey Abramov, and Florentin Wörgötter. Occlusion handling in video segmentation via predictive feedback. In *Computer Vision – ECCV 2012. Workshops and Demonstrations*. 2012.
- [62] Jeremie Papon, Alexey Abramov, Markus Schoeler, and Florentin Wörgötter. Voxel cloud connectivity segmentation - supervoxels for point clouds. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, june 2013.
- [63] Jeremie Papon, Tomas Kulvicius, Eren Erdal Aksoy, and Florentin Wörgötter. Point cloud video object segmentation using a persistent supervoxel world-model. In *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, 2013.
- [64] Sylvain Paris. Edge-preserving smoothing and mean-shift segmentation of video streams. In David Forsyth, Philip Torr, and Andrew Zisserman, editors, *European Conference on Computer Vision (ECCV)*, volume 5303 of *Lecture Notes in Computer Science*, pages 460–473. Springer Berlin / Heidelberg, 2008.
- [65] Karl Pauwels, Norbert Krüger, Markus Lappe, Florentin Wörgötter, and Marc M. Van Hulle. A cortical architecture on parallel hardware for motion processing in real time. *Journal of Vision*, 10(10), 2010. doi: 10.1167/10.10.18.
- [66] J. Piaget and B. Inhelder. *The child's conception of space*. W.W. Norton, New York, 1967.
- [67] S Pinker. *The Language Instinct*. Harper Perennial Modern Classics, New York, 1994.
- [68] Steven Pinker and Stephen M Kosslyn. The representation and manipulation of three-dimensional space in mental images. *Journal of Mental Imagery*, 2(1):69–84, 1978.

- [69] P. Pérez, C. Hue, J. Vermaak, and M. Gangnet. Color-based probabilistic tracking. In Anders Heyden, Gunnar Sparr, Mads Nielsen, and Peter Johansen, editors, *European Conference on Computer Vision (ECCV)*, volume 2350 of *Lecture Notes in Computer Science*, pages 661–675. Springer Berlin / Heidelberg, 2002.
- [70] D.B. Reid. An algorithm for tracking multiple targets. *Automatic Control, IEEE Transactions on*, 24(6):843–854, Dec 1979.
- [71] Xiaofeng Ren and J. Malik. Learning a classification model for segmentation. In *ICCV*, pages 10–17 vol.1, 2003.
- [72] Andreas Richtsfeld, Thomas Morwald, Johann Prankl, Michael Zillich, and Markus Vincze. Segmentation of unknown objects in indoor environments. In *IROS*, pages 4791–4796, 2012. ISBN 978-1-4673-1737-5.
- [73] Jürgen Rossmann, Nils Wantia, Eren Erdal Aksoy, Simon Haller, and Justus Piater. Active learning of manipulation sequences. In *Robotics and Automation (ICRA) 2014. IEEE International Conference on*, 2014.
- [74] R.B. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (fpfh) for 3d registration. In *Robotics and Automation, 2009. ICRA '09. IEEE International Conference on*, pages 3212–3217, May 2009.
- [75] A. Saffari, C. Leistner, J. Santner, M. Godec, and H. Bischof. On-line random forests. In *IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, 2009.
- [76] J. Santner, C. Leistner, A. Saffari, T. Pock, and H. Bischof. Prost: Parallel robust online simple tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [77] D. Schulz, W. Burgard, D. Fox, and A.B. Cremers. Tracking multiple moving targets with a mobile robot using particle filters and statistical data association. In *Robotics and Automation, 2001. Proceedings 2001 ICRA. IEEE International Conference on*, volume 2, pages 1665–1670 vol.2, 2001.
- [78] Shinsuke Shimojo, Gerald H Silverman, and Ken Nakayama. Occlusion and the solution to the aperture problem for motion. *Vision research*, 29(5):619–626, 1989.
- [79] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, pages 746–760, 2012.
- [80] Jan Smisek, Michal Jancosek, and Tomás Pajdla. 3D with kinect. In *ICCV Workshops*, pages 1154–1160, 2011. ISBN 978-1-4673-0062-9.
- [81] Miquel F. Sumsi. *Theory and Algorithms on the Median Graph. Application to Graph-based Classification and Clustering*. PhD thesis, Universitat Autònoma de Barcelona, 2008.
- [82] André Ückermann, Robert Haschke, and Helge Ritter. Real-time 3D segmentation of cluttered scenes for robot grasping. In *Humanoids*, 2012.

- [83] Amelio Vazquez-Reina, Shai Avidan, Hanspeter Pfister, and Eric Miller. Multiple hypothesis video segmentation from superpixel flows. In Kostas Daniilidis, Petros Maragos, and Nikos Paragios, editors, *European Conference on Computer Vision (ECCV)*, volume 6315 of *Lecture Notes in Computer Science*, pages 268–281. Springer Berlin / Heidelberg, 2010.
- [84] Olga Veksler, Yuri Boykov, and Paria Mehrani. Superpixels and supervoxels in an energy optimization framework. In *Proceedings of the 11th European conference on Computer vision: Part V, ECCV 10*, pages 211–224, Berlin, Heidelberg, 2010. Springer-Verlag.
- [85] J. Vermaak, Arnaud Doucet, and P. Perez. Maintaining multimodality through mixture tracking. In *Computer Vision, 2003. Proceedings. Ninth IEEE International Conference on*, pages 1110–1116 vol.2, Oct 2003.
- [86] J. Vermaak, S.J. Godsill, and P. Perez. Monte carlo filtering for multi target tracking and data association. *IEEE Transactions on Aerospace and Electronic Systems*, 41(1):309 – 332, jan. 2005.
- [87] Paul Viola and MichaelJ. Jones. Robust real-time face detection. *IJCV*, 57(2):137–154, 2004. ISSN 0920-5691.
- [88] B.-N. Vo, S. Singh, and A. Doucet. Sequential monte carlo methods for multitarget filtering with random finite sets. *IEEE Transactions on Aerospace and Electronic Systems*, 41(4):1224 – 1245, oct. 2005.
- [89] D. Weikersdorfer, D. Gossow, and M. Beetz. Depth-adaptive superpixels. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 2087–2090, 2012.
- [90] Bo Wu and Ram Nevatia. Detection and tracking of multiple, partially occluded humans by bayesian combination of edgelet based part detectors. *International Journal of Computer Vision*, 75(2):247–266, 2007. ISSN 0920-5691.



# **Appendices**



# A

## The Oculus Vision System

### A.1 MOTIVATION

There is great interest in development of complex vision systems for robotic vision applications. Such research has strict requirements; these systems must operate in real-time, using input from multiple sources, and typically consist of multiple algorithms which work in concert to produce useful output with minimal delay. Consequently, the architecture which binds algorithms and input sources together has become an increasingly important factor. In this Appendix we shall present a vision architecture we developed over the course of the thesis work which uses modular plugins, a novel buffering scheme, and GPU memory optimizations to allow real-time performance of an online vision system, even with complex pipelines and algorithms developed by independent researchers.

A primary concern when developing such complex vision systems lies in how to properly integrate algorithms developed by different researchers, often from multiple institutions. Typically, computer vision researchers develop solutions tailor-made for their particular problem, without concern over the difficulties involved in integrating their particular algorithm into a large system. To ease this integration process, we provided a plugin interface. The plugin system allows independently developed algorithms to communicate with the architecture's central memory management system, interact with the GUI, define their own unique data types, and integrate into systems with plugins developed by other researchers.

Another motivation for developing a vision architecture is the desire to enable the use of complex algorithmic layouts in an online system. In particular, interest in creating loops that

allow high level algorithms (i.e. which come late in the pipeline) to feedback and improve the output of low level vision methods. Traditional online vision pipeline architectures cannot accommodate such loops in an adequate way, as at any given moment each portion of the pipeline is processing data from different instants in time.

Existent vision system architectures also do not support the use of GPUs in a fully integrated way, leading to inefficient use of the device and communication with device memory. The presented method incorporates specially designed GPU data-containers to ensure optimal PCI-bus use through a pre-caching scheme and concurrent memory transfers. In addition to these, extensibility is ensured through an interface which allows user-defined data-container handling, allowing plugin developers to explicitly define how the memory manager shares data between the host and device. In this Appendix we will present an overview of our system, describe a typical system configuration used for robotics, and then give performance figures from a demonstration setup.

## A.2 SYSTEM ARCHITECTURE

Our vision system is a plugin shell which provides an easy-to-use API for interacting with the GUI, memory management system, and visualization components. In order to ensure expandability, such a system must provide straightforward communication and interaction between plugins created independently, while employing strong-typing checks to ensure only valid plugins may be inter-connected. In addition, it must ensure that plugins have the flexibility to define their own methods for visualization. Finally, the system must ensure that each plugin is self-contained, and executes within its own thread (or threads). This is especially important for fast execution on modern processors, where the number of cores can match, or even exceed, the number of plugins one is running.

In the next subsections, we shall describe how our architecture accomplishes these goals while requiring as little computational and communication overhead as possible. Small overhead is especially important in the case of real time video processing, where relatively large images must be processed at fast frame rates.

### A.2.1 EXECUTION FLOW

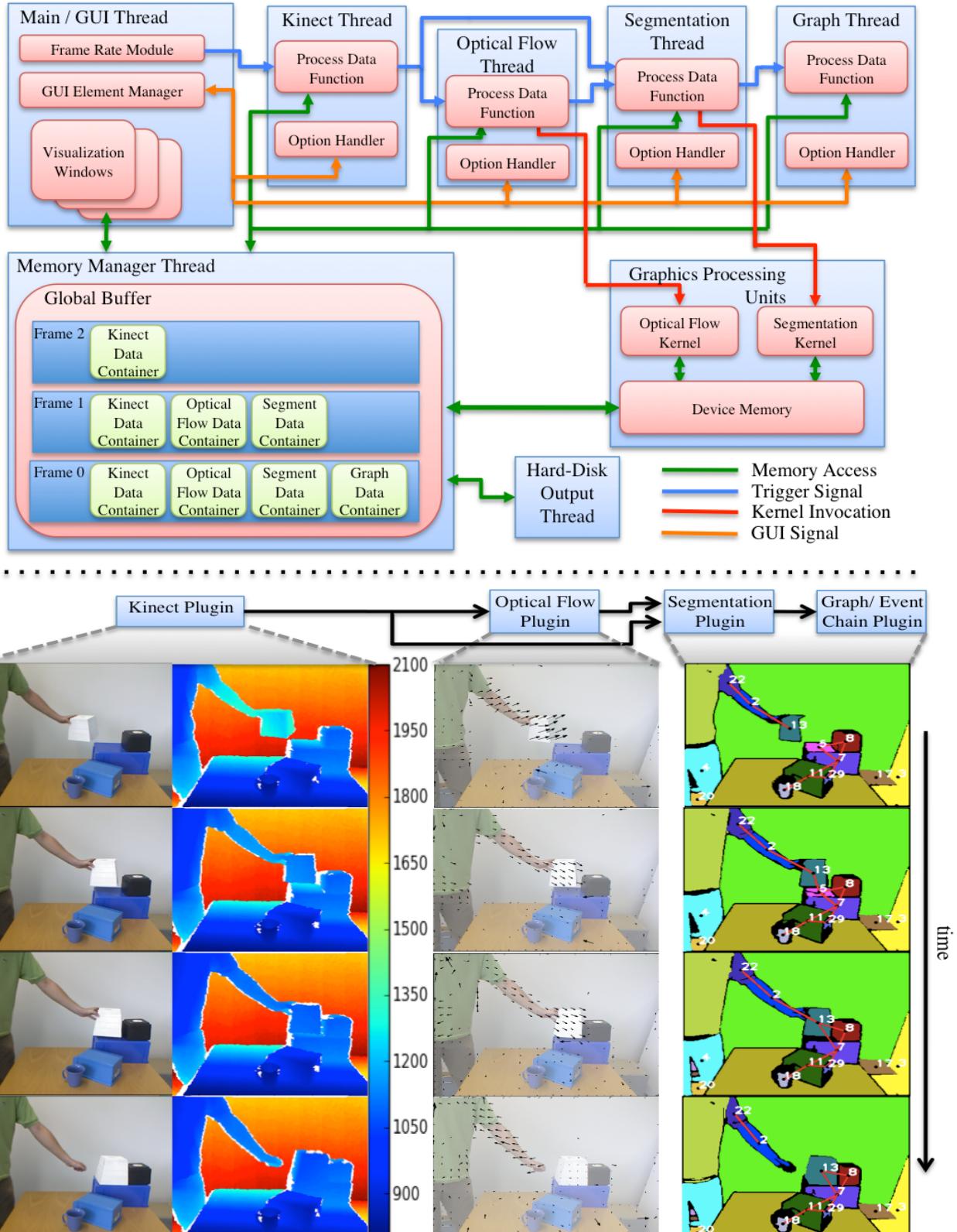
At its core, the architecture provides a shell which consists of a GUI for loading plugins and visualizing data, a system for storing plugin output to file, and a buffering/memory-management system for handling data. This functionality is contained in the *Main Thread* and *Memory Manager Thread* shown in Figure A.2.1. Users build their system by adding plugins, configuring their options via the GUI, and then connecting the plugins to each-other. The user can also save/load a fully configured system as an XML file. Once a vision system has been built, the user can control execution using the frame rate module, which controls the firing rate of the system clock.

As the whole system runs asynchronously in independent threads, the clock trigger acts as the initial starting point for each frame. This means that any source plugins, such as a stereo camera rig or a video file reader plugin, must connect to the frame rate module. As a trigger arrives at each plugin, a triggering signal is sent to the memory manager, telling it to generate a *DataContainer* for the plugin's output. The plugin is then triggered, causing it to execute its processing functionality and generate output, which it stores in the location assigned to it by the memory manager. The plugin then generates another triggering signal, which is connected to both the memory manager and whatever ensuing plugins use the output as their input. When a plugin has multiple inputs, it will loop inside its execution thread, waiting until all inputs for a frame have arrived before executing. This is accomplished by each thread having its own input queue map; it is important to note though, that these queues contain no actual data (and thus minimal overhead), and merely serve as a message passing system. The signaling and triggering system employs the open-source Qt signal & slot architecture. In particular, the system makes use of Qt's ability to queue signals for execution as they arrive at a thread.

#### A.2.2 PLUGIN DEVELOPMENT AND INTERACTION

The functionality of the system is provided primarily via plugins. A plugin consists of a shared library which is loaded dynamically at run-time. The system is based on the low-level Qt plugin API, which facilitates development and ensures compatibility across different platforms. Plugins inherit from a pure abstract interface class which defines a protocol for communicating with the core application. This permits plugins to define input and output types and pass messages to/from the GUI and memory manager.

Developers are required to implement a *processData* function, which receives input and writes to an output *DataContainer*. The developer can optionally create any number of GUI elements (e.g. sliders, buttons) using the interface functions. Plugins specify how many inputs they require, and give the possible types for these inputs. Communication between plugins is accomplished through a standardized data container interface. The core architecture contains commonly used data container implementations, such as *StereoImageContainer*. Plugins may define their own specialized data containers which are loaded at runtime with the plugin. For example, the Segmentation plugin has its own container type *SegmentationData*, which contains a list of labeled segments, metadata about the segments, and labeled images. The standardized data container interface allows for any plugin to refer to a new container class without actual knowledge of the container itself other than the string identifiers of its members (e.g. "Segment Labels"). Correct handling of access to these members is accomplished through dynamic dispatch using the virtual lookup table. This ensures that a plugin written by one researcher can be easily used as input to another's, as long as they know the proper identifiers and underlying formatting of the data.



**Figure A.2.1:** Overview of the system architecture and demonstration system output for four frames. The columns show output from the different components; from left to right, Kinect image and depth (in mm), optical flow, and graphs overlaid on segmentation plugin output. This type of output can be seen live in any number of visualization windows within the GUI.

### A.2.3 VISUALIZATION

During the development and use of a vision system, it is of utmost importance to be able to visualize what is occurring at every stage of the system pipeline. As such, our system allows users to create any number of visualization windows which can select any plugin to display (and which part of the plugin's output to display, e.g. left or right image). If a developer creates their own data container for a plugin, they can define a special visualization callback function as part of this container. The system will automatically detect this callback when the plugin is loaded, and use it for visualizing the plugin's output. Developers can specify multiple methods for visualizing the plugin; the GUI for visualization will allow selection of which to display.

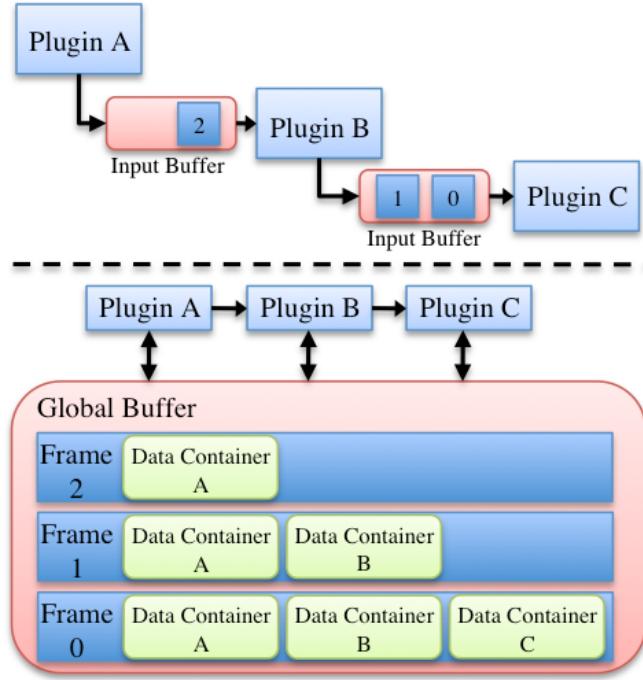
Visualization windows read directly from the global buffer, and as such have a small memory overhead. Additionally, visualization runs in the GUI thread, rather than in any of the plugin threads. If a plugin slows down the system, visualization (and the GUI) will remain responsive, allowing the user to troubleshoot. This also means that visualization that requires computation, such as labeling an image with text or vector graphics, will have a negligible effect on the actual frame throughput of the system. If visualization lags behind the system output, frames are automatically skipped on an interval that allows visualization to maintain synchronization with the rest of the system. This is of particular importance in an online system, such as our real-time robotic application, where visualization lagging behind processing can cause confusion or even errors.

## A.3 MEMORY ARCHITECTURE

The memory management system has been designed to allow distributed development and computing, complex system pipelines incorporating feedback loops, and efficient use of the GPU as a computational resource. The following subsections will describe how these design goals have been achieved by illustrating our *Global Buffer* design and explaining how it manages GPU memory.

### A.3.1 GLOBAL BUFFER

Our global buffer concept was designed to overcome the limitations of standard online vision pipelines. In a standard online pipeline a local buffering scheme is used; each algorithm has an input buffer, where data accumulates while it is waiting to be processed. Such a setup is adequate as long as the pipeline remains unidirectional, but complications arise in using feedback loops. Figure A.3.1 compares a standard pipeline with our global buffer; unlike a typical buffering scheme, our global buffer maintains and manages all memory in a central location (and separate thread). The global buffer is responsible for dynamic allocation of all data containers, maintaining reference counts, and determining when a frame can expire. Since the global buffer is responsible for maintaining memory, plugins use a message passing system to

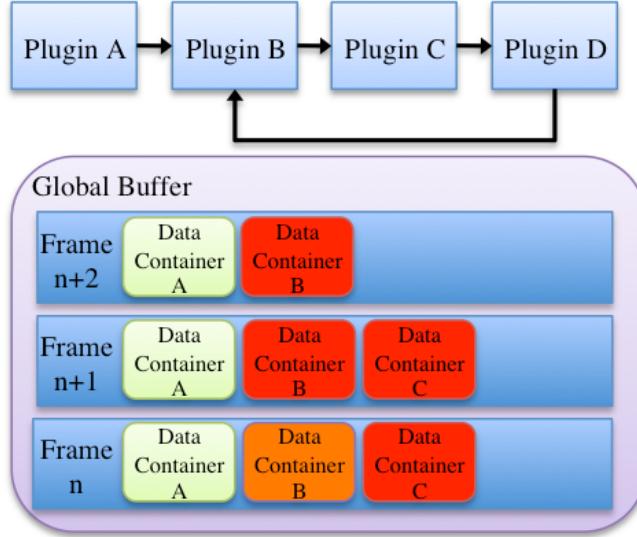


**Figure A.3.1:** A typical buffering scheme (top) and our buffer (bottom).

communicate. Plugins pass messages to each other to notify completion of a new frame, or to trigger a feedback mechanism. They also use the message passing system to request that the global buffer allocate a new data container for their output. When a developer creates a new type of data container, they use a simple interface to pass the global buffer a function pointer for creating an instance of their new data container type.

In order to fully understand the limitations of a standard buffering system, consider, for instance, the system shown at the top of Figure A.3.2. If the feedback mechanism is triggered for frame  $n$ , plugin  $B$  must return to frame  $n$  in order to modify how it was processed. This is not possible in the standard local buffer scheme, as that data was discarded after it was used as input to  $B$ . One possible solution is to maintain another local buffer for each plugin which contains data which has already been processed, but this quickly adds several degrees of complexity. In particular, garbage collection becomes very difficult, and management of these buffers when feedback does occur becomes unnecessarily convoluted.

The global buffer solves this by maintaining data in a more structured way. When a feedback mechanism is triggered for frame  $n$  the triggering plugin ( $D$ ) sends a message to  $B$ , causing it to stop processing what it has scheduled, and revert to frame  $n$ . As frame  $n$  is still easily accessible in the global buffer,  $B$  can simply send a request for the pointer(s) to the input data container(s) it requires. The global buffer is guaranteed to still have the data for frame  $n$ , because  $D$  never produced an output for frame  $n$ , so the global buffer has not marked frame  $n$  as complete. Once  $B$  finishes processing frame  $n$  with its new feedback information, it will overwrite its old output for frame  $n$  (shown in orange) and then simply continue on as it would normally,



**Figure A.3.2:** Feedback using a global buffer

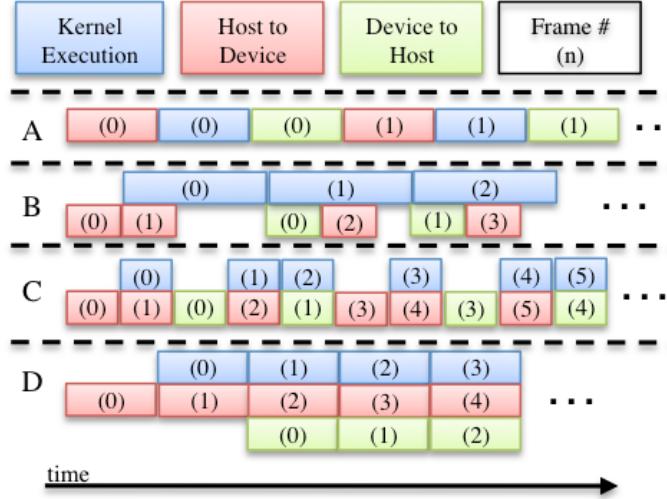
processing frame  $n+1$ . The feedback corrected data will propagate down the pipeline, and any data which is no longer valid (shown in red) will simply be overwritten. Infinite feedback loops are avoided by preventing feedback from occurring more than once per plugin per frame.

### A.3.2 GPU MEMORY HANDLING

While utilizing the massively-parallel GPU as a coprocessor has become increasingly common, how to integrate it effectively into an open vision architecture remains an open question. Particularly vexing is how to integrate it seamlessly into the memory system of such an architecture, as the GPU has separate physical memory, which is entirely distinct in both location and structure from that used by the CPU. Data streaming through the system must be transferred to the GPU for modules which use it, and then transferred back out for visualization and used by modules later in the pipeline.

A naive implementation of this architecture would simply serialize the operations; when a module needs to use the GPU, it copies data to device memory, executes a kernel, and then copies the output back out to host memory. While this is still relatively efficient, it fails to fully take advantage of the pipelined streaming architecture, since the memory transfer bandwidth is idle while the kernel is executing. The architecture uses the streaming CUDA API to utilize this spare bandwidth, allowing it to perform concurrent asynchronous memory transfer and kernel execution.

As shown in Figure A.3.3, we utilize a pre-caching technique, whereby data for frame  $n+1$  is transferred during the execution of frame  $n$ . When the kernel execution time is significantly longer than the transfer time ( $B$ ), memory transfer is completely hidden, even with unidirectional memory. When kernel execution time is comparable to memory transfer time, only



**Figure A.3.3:** Streaming; Concurrent kernel execution

some of the transfer can be hidden (C), unless the hardware supports concurrent data transfers<sup>1</sup> (D).

## A.4 DEMONSTRATION SYSTEM

This section presents a real-time demonstration system, consisting of six plugins. The demonstration system calculates dense disparity using a standard stereo camera setup (rather than Kinect data) in order to show the flexibility of the architecture as well as highlight the speedup achieved via multithreading. Switching from Kinect input to a stereo camera setup is simply a matter of changing connections in the GUI. The pipeline described consists of plugins for reading and rectifying stereo data, calculating optical flow[65], computing disparity[65], segmentation and tracking[4], dense disparity estimation, and semantic graph and event chain generation[8, 10]. This type of a system configuration is used to recognize and learn object manipulation actions in a robotics context.

### A.4.1 IMAGE ACQUISITION

Video is acquired using a Firewire stereo camera rig. Triggering for image acquisition can be controlled using either an external hardware trigger or the architecture's software clock. Rectification is performed on the GPU (there is a separate plugin for calibration using a standard chessboard). Time from triggering to output of a rectified pair of stereo images is around 10ms at 1024x768.

---

<sup>1</sup>Concurrent data transfers are supported under the Fermi architecture. Currently the Fermi Quadro and Tesla series cards have two Direct memory access (DMA) engines, allowing them to perform host-to-device and device-to-host operations simultaneously. The consumer Fermi cards (GTX 4xx, 5xx) only have a single DMA engine, so concurrent transfers are disabled on them.

#### A.4.2 DISPARITY AND OPTICAL FLOW

Optical flow is computed using the GPU implementation [65] of a phase-based algorithm [35]. The algorithm tracks the temporal evolution of equi-phase contours by taking advantage of phase constancy. Differentiation of the equi-phase contours with respect to time yields spatial and temporal phase gradients. Optical flow is then computed by integrating the temporal phase across orientation. Estimates are refined by traversing a Gabor pyramid from coarser to fine levels. The plugin uses the five most recent frames to compute optical flow in the case of online video, but can also use "future" frames when working with recorded movies (this can slightly improve the quality of output flow).

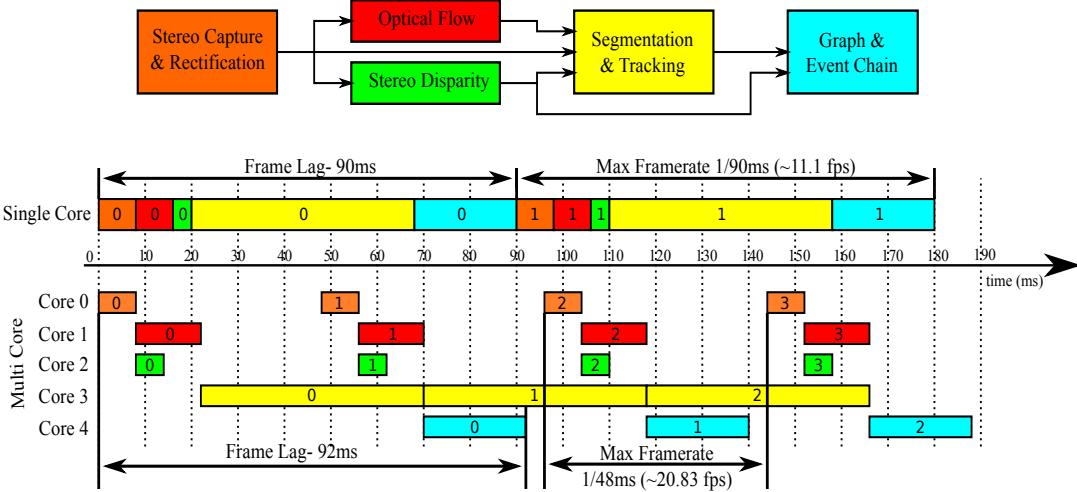
Sparse disparity maps are computed on the GPU using a technique similar to optical flow [65]. Rather than use temporal phase gradients, the disparity algorithm relies on phase differences between stereo-pair rectified images. As with the optical flow algorithm, results are computed using a coarse to fine pyramid scheme.

#### A.4.3 SEGMENTATION AND TRACKING

The segmentation and segment tracking plugin has two roles; first, it partitions the image into labeled regions, as seen in the right-most column of Figure A.2.1, and second, it determines correspondences between frames to maintain consistent labeling. The segmentation algorithm is based on the work of Blatt et al. [17], which applies the Potts model in such a way that superparamagnetic phase regions of aligned spins correspond to a natural partition of the image data. Initial spins are assigned to pixels randomly, and then a Metropolis-Hastings algorithm with annealing [4] is used to iteratively update the spins until an equilibrium state is reached.

The Metropolis algorithm is implemented on the GPU[4], permitting real-time performance. The algorithm itself lends itself to efficient implementation on a GPU, as interactions are only computed locally (8 connected nearest-neighbors). Coupling interactions between pixels are determined using average color vector difference (in the HSV space) of nearest-neighbors. Additionally, when depth data is available, the algorithm prevents interactions between pixels if there is a significant difference in their depth values. This prevents coupling across regions which have similar color but discontinuous depth.

In addition to segmentation, the plugin maintains consistent labels for objects from frame to frame. This is accomplished by transferring spins between frames using output from an optical-flow plugin [4]. As such, only the first frame is actually initialized at random; subsequent frames are initialized using a forward-propagated version of the previous frame's equilibrium spins. This has two advantages; the number of iterations needed to reach equilibrium is greatly reduced since the spin distribution already approximates the final state, and the algorithm naturally tracks objects since spins (and thus labels) are maintained over time.



**Figure A.4.1:** Timing results for demonstration system; plugins are color coded and contain frame numbers. When run in single thread mode, short GPU operations such as optical flow are significantly faster due to reduced overhead; this results in slightly lower (2ms) frame lag. The true benefit of multi-threaded mode is the higher maximum frame-rate that can be achieved.

#### A.4.4 SEMANTIC GRAPHS

The semantic graphs plugin constructs a symbolic 3D description of the scene from the segmentation results and disparity maps. Segments are used to construct undirected and unweighted graphs (seen in the right-most column of Figure A.2.1; nodes are labeled with numbers and red lines are graph edges). Each segment is given a node and edges represent their three dimensional touching relations. Graphs can change by continuous distortions (lengthening or shortening of edges) or, more importantly, through discontinuous changes (nodes or edges can appear or disappear). Such a discontinuous change represents a natural breaking point: All graphs before are topologically identical and so are those after the breaking point. Hence, we can apply an exact graph-matching method [81] at each breaking point and extract the corresponding topological main graphs. The sequence of these main graphs thus represents all structural changes (manipulation primitives) in the scene.

This type of graph representation is then encoded by a semantic event chain (SEC), which is a sequence-table; rows and columns of which represent possible spatial relations between each segment pair and manipulation primitive. This final output can be used to classify manipulations and categorize manipulated objects for use in a robotics or human-computer interaction (HCI) setting[8, 10]. The primary advantage of this method is that actions can be analyzed without models or a-priori representation; the dynamics of an action can be acquired without needing to know the identities of the objects involved.

## A.5 RESULTS AND DISCUSSION

Testing was performed to compare single threaded with multi-threaded operation mode and to detect the impact of visualization on processing speed. Testing was performed on an Intel i7 (3.33Ghz, 8 execution threads) system with an NVIDIA GTX 295 GPU. The demonstration setup depicted at the top of Figure A.4.1 was used for all tests. To determine if visualization had a negative impact, the tests were run with and without a visualization windows for each component, showing live views of their outputs. Timing measurements for plugins are the mean execution time per frame of a 1000 frame (640x480) stereo video sequence (frames of which are shown in Figure A.2.1), averaged over 10 runs. The code for the single and multi-threaded versions is identical with the exception of the movement of plugin objects to separate threads.

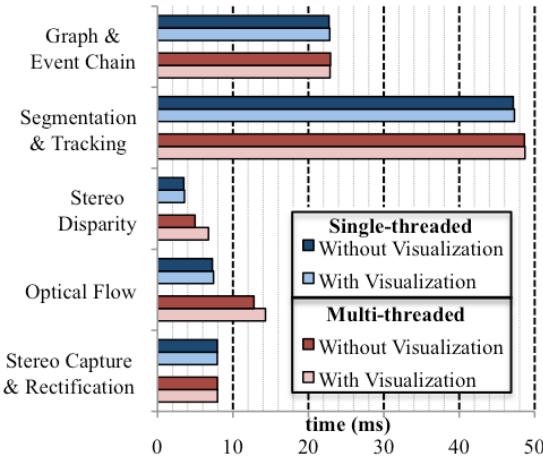
We measure performance by analyzing two key attributes of a pipelined vision real-time vision system. First, in terms of frame lag, that is time from frame acquisition to final output, multi-threaded mode is slightly slower than single-threaded. As shown in Figure A.4.1, this is due to relatively fast plugins which use the GPU (disparity and optical flow in this case). This can be attributed to the static overhead cost incurred by switching between threads while using the CUDA run-time API. The switching is relatively expensive for short GPU operations as it forces the CUDA driver to create and destroy GPU contexts<sup>2</sup> This could be avoided by the addition of an additional GPU; in our demonstration system the driver is forced to change contexts as there are three threads (flow, disparity, segmentation) attempting to use two GPUs. Additionally, the architecture will soon be brought to the newest CUDA release, which allows context sharing between threads. It should also be noted that at higher resolutions multi-threaded mode overtakes single-threaded, as the overhead cost of context switching is outweighed by the gain from computing optical flow and disparity in parallel.

The second measure of performance, throughput, or maximum frame rate, shows a significant speedup in multi-threaded mode, almost doubling from 11.1 (stereo)fps to 20.83. While significant, the speedup is not equal to the number of execution threads used by the demonstration setup (six; one for each plugin and one for the GUI & memory manager). This less-than-optimal gain can be attributed to the fact that the demonstration system had one component, segmentation & tracking, which was significantly slower than the rest. As seen in Figure A.4.1, the entire system throughput is limited by the rate at which the segmentation plugin produces output.

As seen in Figure A.5.1, the addition of visualization components has a small impact on performance. This delay was most noticeable for the shorter components, disparity and optical flow, but never exceeded 2ms. Fortunately, this additional time does not affect throughput in multi-threaded mode, as it is hidden by the length of the longest component. The times with

---

<sup>2</sup>GPU contexts are analogous to CPU processes, and each have their own distinct address space. Each thread may only have one context active at a time, and contexts may not share threads.



**Figure A.5.1:** Visualization has a slight impact on performance, but the effect is negligible in multi-threaded mode where the slight increases in processing time are hidden in the length of the longest component (in this case, segmentation).

visualization were used for Figure A.4.1; clearly shortening the time of any component other than segmentation will have a negligible effect on performance. While the increase does not affect throughput, it has a slight effect on frame lag. Frame lag is less important than throughput for our research, but it should be noted that in certain cases, such as when quick reactions are required, frame lag may be an important performance measure.

## A.6 CONCLUSION

Building a self-contained, efficient, and complete vision system acts as a significant barrier to entry for those wishing to develop and test new vision algorithms. We have presented a modular plugin environment, designed specifically for expandability and parallel architectures, which facilitates rapid distributed development of vision pipelines. Our plugin system allows simple collaboration between organizations, allowing developers to share algorithms easily, and without forcing them to share code. The architecture permits streaming use of the GPU as a coprocessor, efficient visualization of algorithm outputs, and the ability to use complex pipelines involving feedback mechanisms. The system architecture has been released under an open-source GPL license<sup>3</sup>.

---

<sup>3</sup><https://launchpad.net/oculus>

# B

## Sequential Bayesian Estimation

Sequential Bayesian estimation refers to a class of approaches for estimating a varying unknown probability density function from a time series of noisy observations. These approaches use a state space representation, in which a state vector  $\mathbf{x}_t$  describes the hidden state of a dynamic system. The goal is to estimate the posterior distribution of the state given all prior observations  $\mathbf{z}$ , i.e.,  $p(\mathbf{x}_t | \mathbf{z}_{1:t})$ . This is accomplished using a two step recursion which first generates a hypothesis of the current state conditioned on the previous state and then performs a Bayes update using the new observation. These steps are known as the prediction and filtering steps, respectively.

The prediction step estimates the current distribution given all prior observations, or

$$p(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}. \quad (\text{B.1})$$

This requires the specification of a stochastic *dynamic model* to characterize the state transition density  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ :

$$\mathbf{x}_t = f_t(\mathbf{x}_{t-1}, \mathbf{v}_t), \quad (\text{B.2})$$

where  $\mathbf{v}_t$  is the process noise. The dynamic model takes advantage of knowledge of the system to generate reliable predictions of how the state evolves independent of observations.

The filtering step uses Bayes rule to update the predicted density by conditioning it on the new observation  $\mathbf{z}_t$ :

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}) = \frac{p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{z}_{1:t-1})}{p(\mathbf{z}_t | \mathbf{z}_{1:t-1})}. \quad (\text{B.3})$$

This requires the specification of a *measurement model* to characterize the observation density  $p(\mathbf{z}_t | \mathbf{x}_t)$ :

$$\mathbf{z}_t = h_t(\mathbf{x}_t, \mathbf{w}_t), \quad (\text{B.4})$$

where  $\mathbf{w}_t$  is the measurement noise. The marginal likelihood  $p(\mathbf{z}_t | \mathbf{z}_{1:t-1})$  is constant relative to the state, and is generally ignored in practice and replaced with a simple normalizing factor.

Once the filtered, or posterior distribution is determined, an estimate of the state can be made using a variety of techniques (e.g., MAP, mean-shift).

## B.1 PARTICLE FILTERS

Unfortunately, except for in special cases (such as the linear Gaussian case with the Kalman filter) determining an exact solution for the posterior distribution is not feasible. As such, Particle Filter techniques were developed to approximate the posterior distribution. They use sequential Monte Carlo to directly implement the Bayesian recursion equations on a set of samples. The most common Particle Filter algorithm is Sequential Importance Sampling (SIS) recursively updates a set of  $N$  samples (particles) from the previous time step  $\{\mathbf{x}_{t-1}^j, w_{t-1}^j\}$  in a two-step procedure:

1. **Predict:** Apply the dynamic model to find an estimate of the new state for each particle,  $\tilde{\mathbf{x}}_t^{1..N}$ . That is, draw samples from the state transition prior distribution  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ .
2. **Update:** Evaluate the weight for each particle using the observation density:  $\tilde{w}_t^j = p(\mathbf{z}_t | \tilde{\mathbf{x}}_t)$  and then normalize.

The set of weighted particles  $\{\mathbf{x}_t^j, w_t^j\}$  then approximates the posterior distribution, and an overall state estimate can be found using any appropriate method.

### B.1.1 RESAMPLING

An important issue with SIS is that for any finite number of particles the weights will tend to degenerate to the trivial set where all particles have weight zero except for one. This results in the observations having no effect on the particle trajectories, meaning the filter amounts to a random walk using the dynamic model. To avoid this problem, a resampling step was added [36] which generates a new particle set by sampling from the existing particle set. The simplest way of doing this is to simply sample from the multinomial distribution of the particle weights and then set all particle weights to  $1/N$ . While this *multinomial resampling* can be effective if employed judiciously, it can also lead to other problems, namely an increasing variance of the posterior distribution. To overcome this a variety of low-variance resampling techniques have been developed; we refer the reader to [30] for a description of different approaches.