

16.485: VNAV - Visual Navigation for Autonomous Vehicles

Luca Carlone



Lecture 27: Advanced Topics: Dense 3D Reconstruction



Big Picture



Big Picture



Today

- **Dense Reconstruction**
 - 3D representations
 - (Some) Multi-view Stereo
 - Depth fusion
- Final thoughts

Multi-View Stereo: A Tutorial 2015

Yasutaka Furukawa

Washington University in St. Louis furukawa@wustl.edu

> Carlos Hernández Google Inc. carloshernandez@google.com

ElasticFusion: Dense SLAM Without A Pose Graph

Thomas Whelan*, Stefan Leutenegger*, Renato F. Salas-Moreno[†], Ben Glocker[†] and Andrew J. Davison* *Dyson Robotics Laboratory at Imperial College, Department of Computing, Imperial College London, UK [†]Department of Computing, Imperial College London, UK

{t.whelan,s.leutenegger,r.salas-moreno10,b.glocker,a.davison}@imperial.ac.uk

2016

KinectFusion: Real-Time Dense Surface Mapping and Tracking*

Richard A. Newcombe Imperial College London

Andrew J. Davison Imperial College London

Shahram Izadi Otmar Hilliges Microsoft Research Microsoft Research

Pushmeet Kohli Microsoft Research

Jamie Shotton

Microsoft Research

David Molyneaux David Kim Microsoft Research Microsoft Research Lancaster University Newcastle University Steve Hodges Andrew Fitzgibbon Microsoft Research Microsoft Research





Figure 1: Example output from our system, generated in real-time with a handheld Kinect depth camera and no other sensing infrastructure. Normal maps (colour) and Phong-shaded renderings (greyscale) from our dense reconstruction system are shown. On the left for comparison is an example of the live, incomplete, and noisy data from the Kinect sensor (used as input to our system).

Voxblox: Incremental 3D Euclidean Signed Distance Fields for On-Board MAV Planning

Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto Autonomous Systems Lab, ETH Zürich

Abstract-Micro Aerial Vehicles (MAVs) that operate in unstructured, unexplored environments require fast and flexible local planning, which can replan when new parts of the map are explored. Trajectory optimization methods fulfill these needs, but require obstacle distance information, which can be given by Euclidean Signed Distance Fields (ESDFs).

We propose a method to incrementally build ESDFs from Truncated Signed Distance Fields (TSDFs), a common implicit surface representation used in computer graphics and vision. TSDFs are fast to build and smooth out sensor noise over many observations, and are designed to produce surface meshes. Meshes allow human operators to get a better assessment of the robot's environment, and set high-level mission goals.



Point Clouds





Point Clouds



Map representation	3D Topology?	Lightweight?	Filters Noise/ Outliers?	Semantics?	Generality	
Point Clouds	×	√/X No, if Dense	X	✓/X No, if Sparse	\checkmark	



Map representation	3D Topology?	Lightweight?	Filters Noise/ Outliers?	Semantics?	Generality
Point Clouds	×	√/X No, if Dense	X	✓/X No, if Sparse	\checkmark
Geometric primitives	×	\checkmark	\checkmark	√/X No, if Sparse	×

Object-based Maps



Volumetric Methods: Voxels/Octrees



Map representation	3D Topology?	Lightweight?	Filters Noise/Outliers?	Semantics?	Generality
Point Clouds	×	✓/X No, if Dense	X	✓/X No, if Sparse	√
primitives & objects	×	√	√	✓/X No, if Sparse	×
Voxels	\checkmark	✓/X No, if small voxel	\checkmark	✓/X No, if large	\checkmark

Meshes





Map representation	3D Topology ?	Lightweight?	Filters Noise/ Outliers?	Semantics?	Generalit y
Point Clouds	×	✓/X No, if Dense	X	✓/X No, if Sparse	\checkmark
primitives & objects	×	\checkmark	\checkmark	✓/X No, if Sparse	×
Voxels	\checkmark	✓/X No, if small voxel	\checkmark	✓/X No, if large voxel	\checkmark
3D Mesh	\checkmark	\checkmark	×	\checkmark	\checkmark

Today

- Dense Reconstruction
 - 3D representations —
 - (Some) Multi-view Stereo
 - Depth fusion
- Final thoughts

Multi-View Stereo: A Tutorial 2015

Yasutaka Furukawa

Washington University in St. Louis furukawa@wustl.edu

> Carlos Hernández Google Inc.

carloshernandez@google.com

ElasticFusion: Dense SLAM Without A Pose Graph

Thomas Whelan*, Stefan Leutenegger*, Renato F. Salas-Moreno[†], Ben Glocker[†] and Andrew J. Davison* *Dyson Robotics Laboratory at Imperial College, Department of Computing, Imperial College London, UK [†]Department of Computing, Imperial College London, UK

{t.whelan,s.leutenegger,r.salas-moreno10,b.glocker,a.davison}@imperial.ac.uk

2016

KinectFusion: Real-Time Dense Surface Mapping and Tracking*

Richard A. Newcombe Imperial College London

Andrew J. Davison Imperial College London

Shahram Izadi Otmar Hilliges Microsoft Research Microsoft Research

Pushmeet Kohli Microsoft Research

Jamie Shotton Microsoft Research

Microsoft Research Lancaster University Steve Hodges Microsoft Research

David Molyneaux Microsoft Research Newcastle University Andrew Fitzgibbon Microsoft Research

David Kim





Figure 1: Example output from our system, generated in real-time with a handheld Kinect depth camera and no other sensing infrastructure. Normal maps (colour) and Phong-shaded renderings (greyscale) from our dense reconstruction system are shown. On the left for comparison is an example of the live, incomplete, and noisy data from the Kinect sensor (used as input to our system).

Voxblox: Incremental 3D Euclidean Signed Distance Fields for On-Board MAV Planning

Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto Autonomous Systems Lab, ETH Zürich

Abstract-Micro Aerial Vehicles (MAVs) that operate in unstructured, unexplored environments require fast and flexible local planning, which can replan when new parts of the map are explored. Trajectory optimization methods fulfill these needs, but require obstacle distance information, which can be given by Euclidean Signed Distance Fields (ESDFs).

We propose a method to incrementally build ESDFs from Truncated Signed Distance Fields (TSDFs), a common implicit surface representation used in computer graphics and vision. TSDFs are fast to build and smooth out sensor noise over many observations, and are designed to produce surface meshes. Meshes allow human operators to get a better assessment of the robot's environment, and set high-level mission goals.



Multi-view Stereo

[courtesy of N. Snavely]

From previous lectures: we know how to use SLAM to get a good estimate of the poses of the cameras



Multi-view stereo

Multi-view Stereo

Towards Internet-scale Multi-view Stereo

CVPR 2010

Yasutaka Furukawa¹ Brian Curless² Steven M. Seitz^{1,2} Richard Szeliski³

> Google Inc.¹ University of Washington² Microsoft Research³

The Visual Turing Test for Scene Reconstruction Supplementary Video

> Qi Shan⁺ Riley Adams⁺ Brian Curless⁺ Yasutaka Furukawa^{*} Steve Seitz^{+*}

⁺University of Washington ^{*}Google

3DV 2013

Multi-view Stereo



Figure 2. Definition of a patch (left) and of the images associated with it (right). See text for the details.

Estimate normal and center of patch to maximize **photometric consistency**:

$$C_{ij}(p) = \rho(I_i(\Omega(\pi_i(p))), I_j(\Omega(\pi_j(p)))))$$

$$Matching \qquad Image \qquad Rectangular \qquad 3D point \\ Score \qquad Intensity \qquad Patch \qquad Projection To camera \\ To c$$

Example of matching score:

$$1 - \sum_{x,y} |W_1(x,y) - W_2(x,y)|^2$$

[Furukawa and Ponce, "Accurate, Dense, and Robust Multi-View Stereopsis", 2007]

Multi-view Stereo

Enforcing regularity: Markov Random Fields Find depth k_p of point "p" such that point is photo-consistent <u>and</u> <u>depth changes smoothly</u>.

$$E(\{k_p\}) = \sum_{p} \Phi(k_p) + \sum_{(p,q) \in \mathcal{N}} \Psi(k_p, k_q)$$

Unary potentials
(similar to previous slides)
$$\Phi(k_p = d) = \min(\tau_u, 1 - \mathcal{C}(p, d))$$

Pairwise potentials
$$\Psi(k_p = d_1, k_q = d_2) = \min(\tau_p, |d_1 - d_2|)$$

Depth is typically discretized before solving..



How Accurate is Multi-view Stereo?





Space Carving Results: African Violet





Input Image (1 of 45) Reconstruction



Reconstruction

Reconstruction

Comparison: 90% of points within 0.128 m of laser scan (building height 51m)

Space Carving Results: Hand



Views of Reconstruction

M. Goesele, N. Snavely, B. Curless, H. Hoppe, S. Seitz, Multi-View Stereo for Community Photo Collections, ICCV 2007

Many methods: volumetric stereo, space carving, Shape from silhouettes, carved visual hull

Today

- Dense Reconstruction
 - 3D representations —
 - (Some) Multi-view Stereo
 - Depth fusion
- Final thoughts

Multi-View Stereo: A Tutorial 2015

Yasutaka Furukawa

Washington University in St. Louis furukawa@wustl.edu

Carlos Hernández

Google Inc. carloshernandez@google.com

ElasticFusion: Dense SLAM Without A Pose Graph

Thomas Whelan*, Stefan Leutenegger*, Renato F. Salas-Moreno[†], Ben Glocker[†] and Andrew J. Davison* *Dyson Robotics Laboratory at Imperial College, Department of Computing, Imperial College London, UK [†]Department of Computing, Imperial College London, UK

{t.whelan,s.leutenegger,r.salas-moreno10,b.glocker,a.davison}@imperial.ac.uk

2016

KinectFusion: Real-Time Dense Surface Mapping and Tracking*

Richard A. Newcombe Imperial College London

Andrew J. Davison Imperial College London

Shahram Izadi Otmar Hilliges Microsoft Research Microsoft Research

Pushmeet Kohli Microsoft Research

Jamie Shotton Microsoft Research

David Molyneaux David Kim Microsoft Research Microsoft Research Lancaster University Newcastle University Steve Hodges Andrew Fitzgibbon Microsoft Research Microsoft Research





Figure 1: Example output from our system, generated in real-time with a handheld Kinect depth camera and no other sensing infrastructure. Normal maps (colour) and Phong-shaded renderings (greyscale) from our dense reconstruction system are shown. On the left for comparison is an example of the live, incomplete, and noisy data from the Kinect sensor (used as input to our system).

Voxblox: Incremental 3D Euclidean Signed Distance Fields for On-Board MAV Planning

Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto Autonomous Systems Lab, ETH Zürich

Abstract-Micro Aerial Vehicles (MAVs) that operate in unstructured, unexplored environments require fast and flexible local planning, which can replan when new parts of the map are explored. Trajectory optimization methods fulfill these needs, but require obstacle distance information, which can be given by Euclidean Signed Distance Fields (ESDFs).

We propose a method to incrementally build ESDFs from Truncated Signed Distance Fields (TSDFs), a common implicit surface representation used in computer graphics and vision. TSDFs are fast to build and smooth out sensor noise over many observations, and are designed to produce surface meshes. Meshes allow human operators to get a better assessment of the robot's environment, and set high-level mission goals.





ElasticFusion: Dense SLAM Without A Pose Graph

Thomas Whelan, Stefan Leutenegger, Renato Salas-Moreno, Ben Glocker, Andrew Davison

Imperial College London

note: based on RGB-D (contrarily to multi-view stereo)

A Gentle Start: 2D Occupancy Grid Maps



- discretize the environment into cells
- Each cell holds real number [0,1], representing the probability of the cell being occupied



A Gentle Start: 2D Occupancy Grid Maps

$$p(m \mid z_{1:t}, x_{1:t}) \qquad p(\mathbf{m}_i \mid z_{1:t}, x_{1:t}) \qquad \underset{\text{being occupied}}{\text{Probability of cell being occupied}}$$
Bayes rule (omitting "x" for simplicity):
$$p(m_i \mid z_{1:t+1}) = \frac{p(z_{t+1} \mid m_i)p(m_i \mid z_{1:t})}{p(m_i)} \qquad \underbrace{p(m_i \mid z_{1:t})}_{\text{Prior}}$$
Log-odd representation is typically used to avoid numerical instabilities

$$\frac{p(\mathbf{m}_i \mid z_{1:t}, x_{1:t})}{-p(\mathbf{m}_i \mid z_{1:t}, x_{1:t})} \quad \blacklozenge \quad l_{t,i} = \log \frac{p(\mathbf{m}_i \mid z_{1:t}, x_{1:t})}{1 - p(\mathbf{m}_i \mid z_{1:t}, x_{1:t})}$$

Truncated Signed Distance Function (SDF)

- Store distance to nearest obstacle (with sign)
- Only update around obstacle itself

(implicit surface model)

Update rule:

$$d(\mathbf{x}, \mathbf{p}, \mathbf{s}) = \|\mathbf{p} - \mathbf{x}\| \operatorname{sign} \left((\mathbf{p} - \mathbf{x}) \bullet (\mathbf{p} - \mathbf{s}) \right) (1)$$

$$w_{\operatorname{const}}(\mathbf{x}, \mathbf{p}) = 1$$
(2)

$$D_{i+1}(\mathbf{x}, \mathbf{p}) = \frac{W_i(\mathbf{x})D_i(\mathbf{x}) + w(\mathbf{x}, \mathbf{p})d(\mathbf{x}, \mathbf{p})}{W_i(\mathbf{x}) + w(\mathbf{x}, \mathbf{p})}$$
(3)

$$W_{i+1}(\mathbf{x}, \mathbf{p}) = \min \left(W_i(\mathbf{x}) + w(\mathbf{x}, \mathbf{p}), W_{\max} \right)$$
(4)



[Curless and Levoy, "A Volumetric Method for Building Complex Models from Range Images", 2007]

Kinect Fusion (2011)

SIGGRAPH Talks 2011 **KinectFusion:** Real-Time Dynamic 3D Surface Reconstruction and Interaction

Shahram Izadi 1, Richard Newcombe 2, David Kim 1,3, Otmar Hilliges 1, David Molyneaux 1,4, Pushmeet Kohli 1, Jamie Shotton 1, Steve Hodges 1, Dustin Freeman 5, Andrew Davison 2, Andrew Fitzgibbon 1

1 Microsoft Research Cambridge 2 Imperial College London 3 Newcastle University 4 Lancaster University 5 University of Toronto

GPU, memory ...

Kintinuous (2013)



GPU, bounded memory ...

VoxBlox (2017)

Voxblox: Building 3D Signed Distance Fields for Planning Helen Oleynikova, Zachary Taylor, Marius Fehr, Juan Nieto, and Roland Siegwart



CPU, memory

From Voxels to Meshes

Marching cubes



https://www.youtube.com/watch?v=B_xk71YopsA

New Representations: Neural Implicit Surfaces

Use neural networks to define an **implicit** surface representation





Signed distance To obstacles





DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation, Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove, CVPR 2019.

(c)

Neural Radiance Fields or NeRF

Use neural networks to regress not a signed distance function, but **color -> rendering**

NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

Ben Mildenhall* UC Berkeley Pratul P. Srinivasan* UC Berkeley Matthew Tancik* UC Berkeley Jonathan T. Barron Google Research Ravi Ramamoorthi UC San Diego Ren Ng UC Berkeley

* Denotes Equal Contribution









Neural Volume Rendering

Frank Dellaert

Publications

Teaching Talks

alks Blog Posts



Frank Dellaert

Professor, Robotics & Computer Vision

- Atlanta, GA
- Georgia Tech
- 🖂 Email
- Twitter
- in LinkedIn
- Github
- YouTube

🞓 Google Scholar

ORCID

NeRF Explosion 2020

21 minute read

Published: December 16, 2020

▶ 0:00				*2	8	:

The result that got me hooked on wanting to know everything about NeRF :-).

Today

- Dense Reconstruction
 - 3D representations —
 - (Some) Multi-view Stereo
 - Depth fusion
- Final thoughts

Multi-View Stereo: A Tutorial 2015

Yasutaka Furukawa

Washington University in St. Louis furukawa@wustl.edu

> Carlos Hernández Google Inc. carloshernandez@google.com

ElasticFusion: Dense SLAM Without A Pose Graph

Thomas Whelan*, Stefan Leutenegger*, Renato F. Salas-Moreno[†], Ben Glocker[†] and Andrew J. Davison* *Dyson Robotics Laboratory at Imperial College, Department of Computing, Imperial College London, UK [†]Department of Computing, Imperial College London, UK

{t.whelan,s.leutenegger,r.salas-moreno10,b.glocker,a.davison}@imperial.ac.uk

2016

KinectFusion: Real-Time Dense Surface Mapping and Tracking*

Richard A. Newcombe Imperial College London

Andrew J. Davison Imperial College London

Shahram Izadi Otmar Hilliges Microsoft Research Microsoft Research

Pushmeet Kohli Microsoft Research

Jamie Shotton Microsoft Research

Microsoft Research Lancaster University Steve Hodges Microsoft Research

David Molyneaux Microsoft Research Newcastle University Andrew Fitzgibbon Microsoft Research

David Kim





Figure 1: Example output from our system, generated in real-time with a handheld Kinect depth camera and no other sensing infrastructure. Normal maps (colour) and Phong-shaded renderings (greyscale) from our dense reconstruction system are shown. On the left for comparison is an example of the live, incomplete, and noisy data from the Kinect sensor (used as input to our system).

Voxblox: Incremental 3D Euclidean Signed Distance Fields for On-Board MAV Planning

Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto Autonomous Systems Lab, ETH Zürich

Abstract-Micro Aerial Vehicles (MAVs) that operate in unstructured, unexplored environments require fast and flexible local planning, which can replan when new parts of the map are explored. Trajectory optimization methods fulfill these needs, but require obstacle distance information, which can be given by Euclidean Signed Distance Fields (ESDFs).

We propose a method to incrementally build ESDFs from Truncated Signed Distance Fields (TSDFs), a common implicit surface representation used in computer graphics and vision. TSDFs are fast to build and smooth out sensor noise over many observations, and are designed to produce surface meshes. Meshes allow human operators to get a better assessment of the robot's environment, and set high-level mission goals.



Robot Perception or Computer Vision?

Computer vision

.. "a day on a cluster with 500 compute cores"





50-100ms latency, embedded, incremental

No longer a dichotomy for many vision applications!

Robot Perception or Computer Vision?



Unordered Vs Sequential



Robot Perception or Computer Vision?

Perception serves action (and vice-versa!)

