# A Brief Introduction to 3D Capture Technology

Mingsong Dou Google 12/05/2017

# My experiences with 3D Capture



Autostereoscopic Display + 3D Data Acquisition

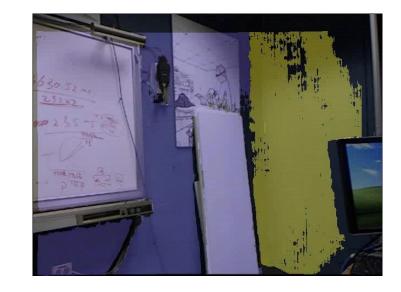
It starts with the Telepresence Project...

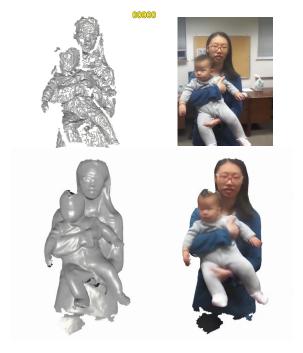
# My experiences with 3D Capture

Gaze-correction for room-sized telepresence, IEEE VR 2012



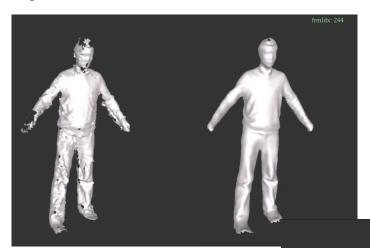
Room scanning with bundle adjustment of points and planes, ACCV 2012



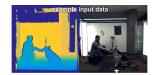


Non-rigid Surface Scanning with dense nonrigid bundle adjustment, CVPR 2015

#### Scanning and Tracking Dynamic Objects, ISMAR 2013









real-time multi-view reconstruction



Holoportation UIST'16



Room-sized Dynamic Scence Capture, IEEE VR 2014

Motion2Fusion SIGGRAPH Asia, 2017



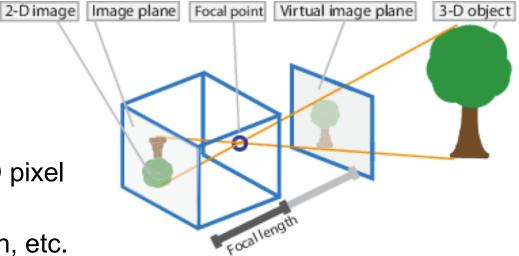


# Outline

- 3D Capture Sensors and Depth Estimation
  - Stereo
  - Structured Light
  - Time-of-Flight
  - Multive-view capture
- World Reconstruction
  - SLAM
  - Kinect Fusion
- People Reconstruction
  - Parametric Tracking
  - Non-Rigid Tracking and Fusion
- Applications in Mixed/Virtual Reality
  - Holoportation

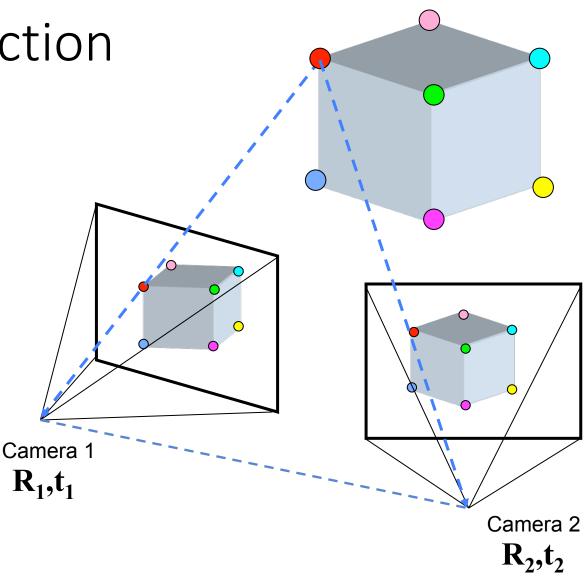
# Camera Pinhole Model

- Camera model:
  - Map a 3D point X=[x; y; z] in the world to a 2D pixel position on the image x=[u; v]
  - Intrinsics: focus length, pixel size, lens distortion, etc.
     Usually represented as a 3x3 matrix *K*.
  - Extrinsics: camera position and orientation, represented by a rotation matrix *R* and translation vector *t*.

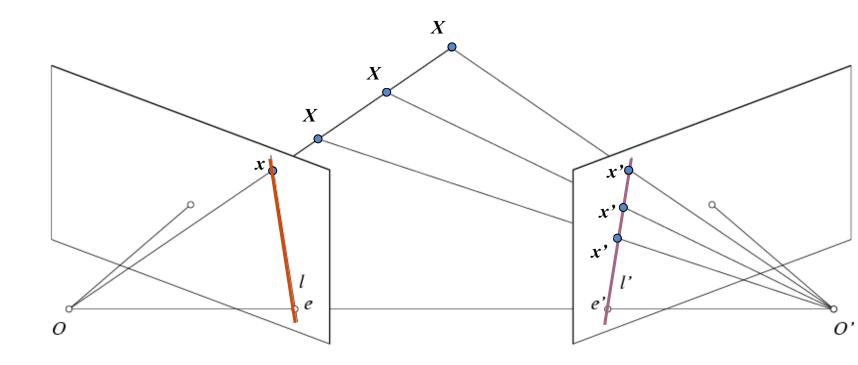


# 3D Capture/Reconstruction

- Reverse Rendering Problem
  - From 2d image point *x* to the corresponding 3D world point *X*
  - Triangulation

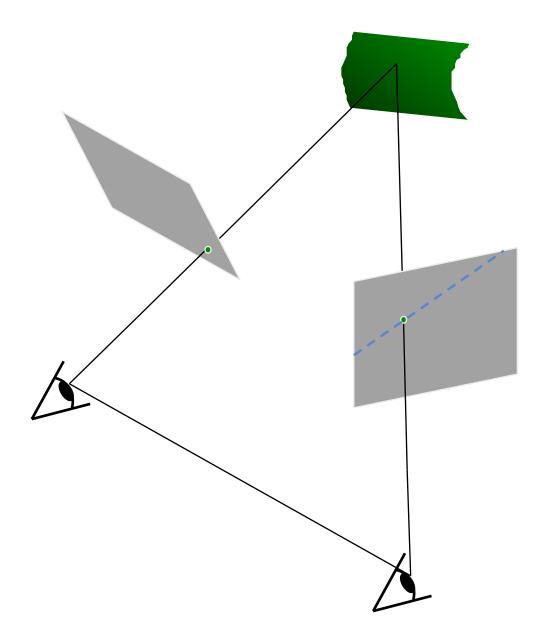


- 1D search problem:
  - search along the epipolar line

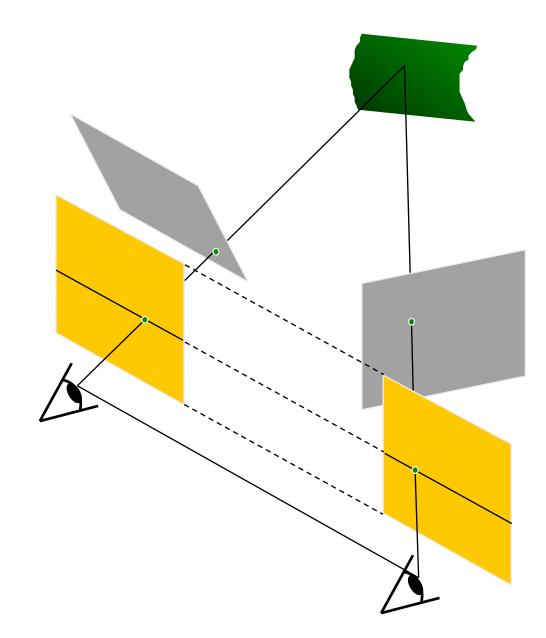


**Epipolar Constraints:** Potential matches for **x** have to lie on the corresponding epipolar line **I**'.

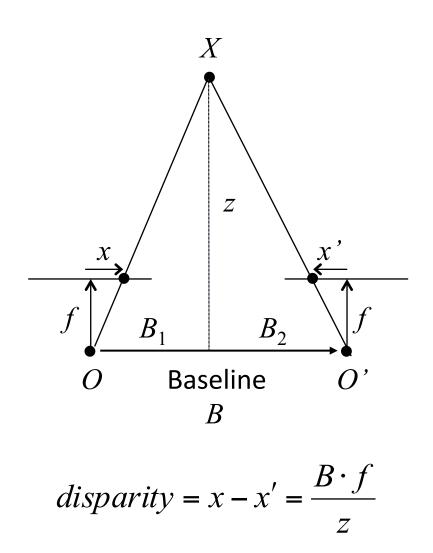
- 1D search problem:
  - search along the epipolar line



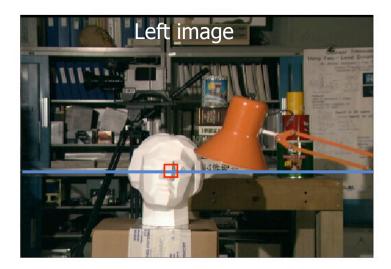
- 1D search problem:
  - search along the epipolar line
- A trick to boost the performance
  - Image rectification

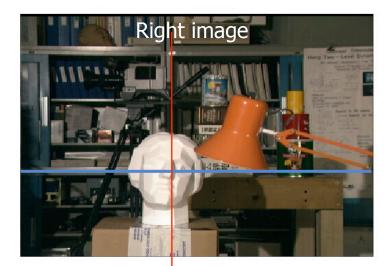


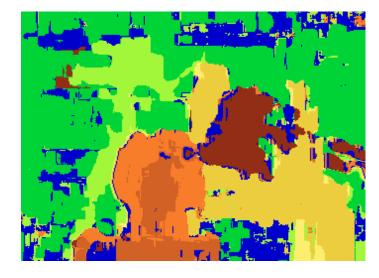
- 1D search problem:
  - search along the epipolar line
- A trick to boost the performance
  - Image rectification
- Disparity map
  - Disparity is inversely proportional to depth

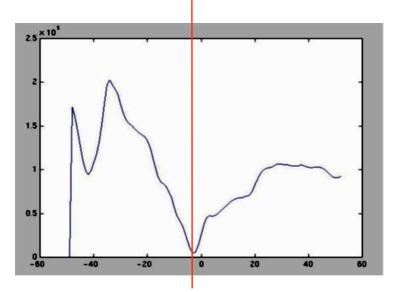


- 1D search problem:
  - search along the epipolar line
- A trick to boost the performance
  - Image rectification
- Disparity map
  - Disparity is inversely proportional to depth

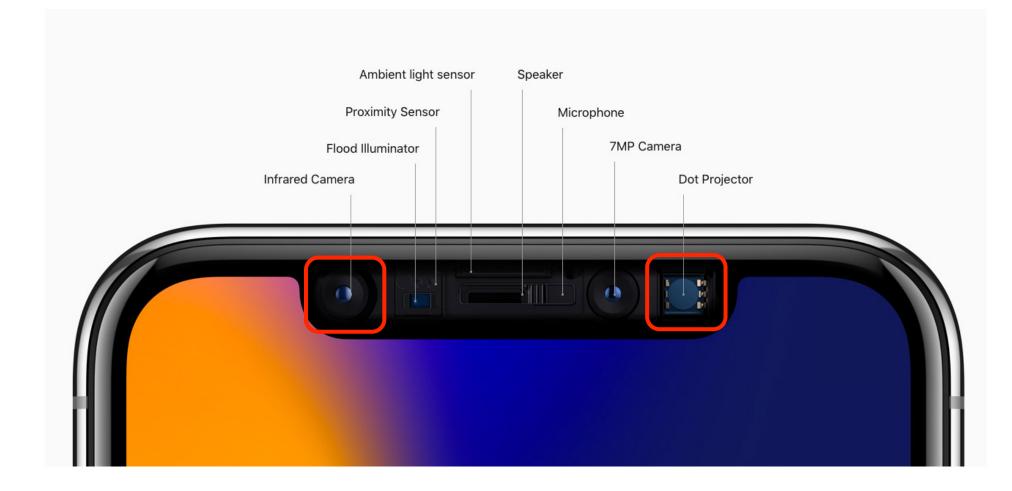


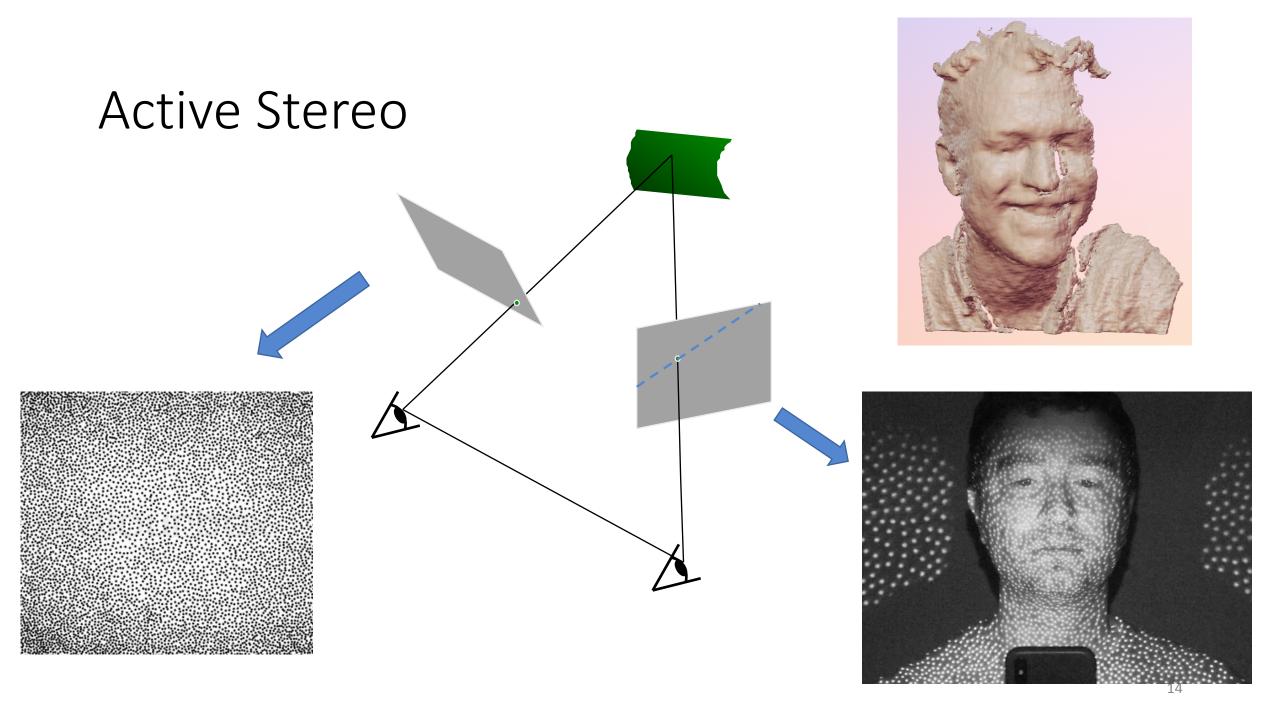






#### Active Stereo



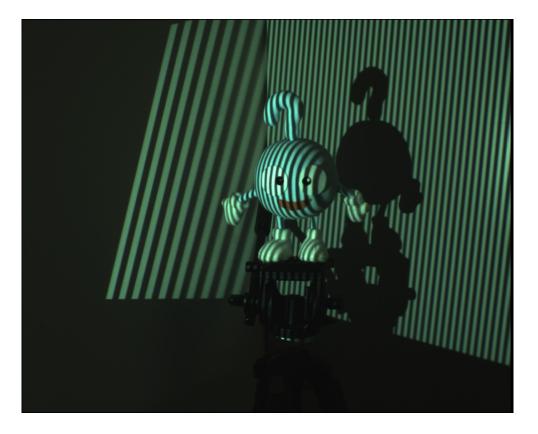


#### Active Stereo

- Other commercial sensors
  - Microsoft Kinect, Intel RealSense, ...
- IR: doesn't work well at outdoor enviroment

#### Other 3D capture techniques

• Structured Light



# Other 3D capture techniques

- Structured Light
- Time-of-flight



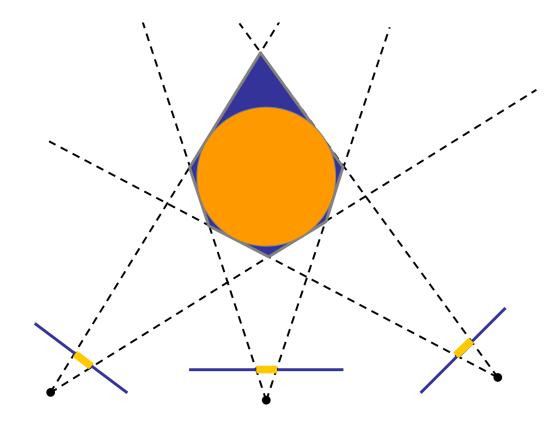
Microsoft Kinect V2

# Other 3D capture techniques

- Structured Light
- Time-of-flight
- LIDAR



#### Shape-from-silhouette



Intersection of foreground cones

# Shape-from-silhouette



Offline, Controlled environment Geometry quality is low, but sharp edges

# Shape-from-silhouette

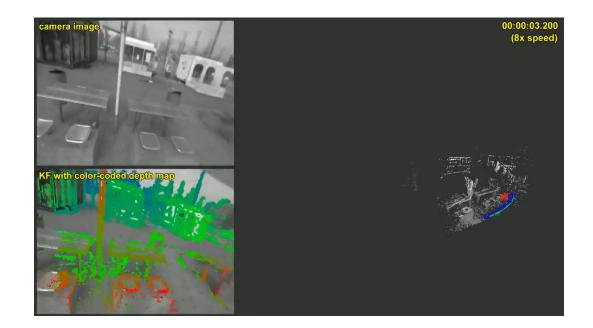


Offline, Controlled environment Geometry quality is low, but sharp edges

# Microsoft Free Viewpoint Video (H-Cap)

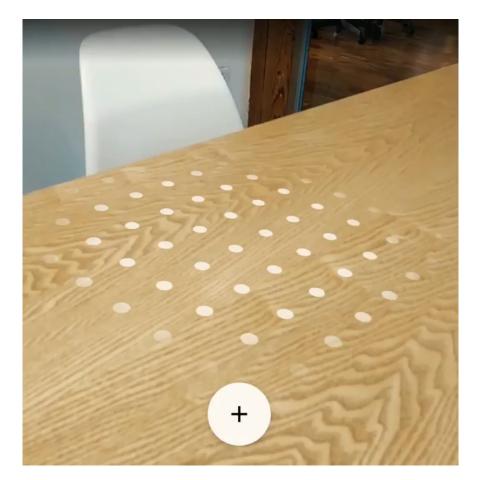


- Simultaneous localization and mapping (SLAM)
  - One moving camera



"LSD-SLAM: Large-Scale Direct Monocular SLAM", J. Engel, T. Schöps, D. Cremers, *ECCV*, 2014

- Simultaneous localization and mapping (SLAM)
  - One moving camera
- Google ARCore/Apple ARKit
  - mobile camera + IMU



Google ARCore

- Simultaneous localization and mapping (SLAM)
  - One moving camera
- Google ARCore/Apple ARKit
  mobile camera + IMU
- Bundle Adjustment
  - a technique to simultaneously optimize both geometry and camera poses



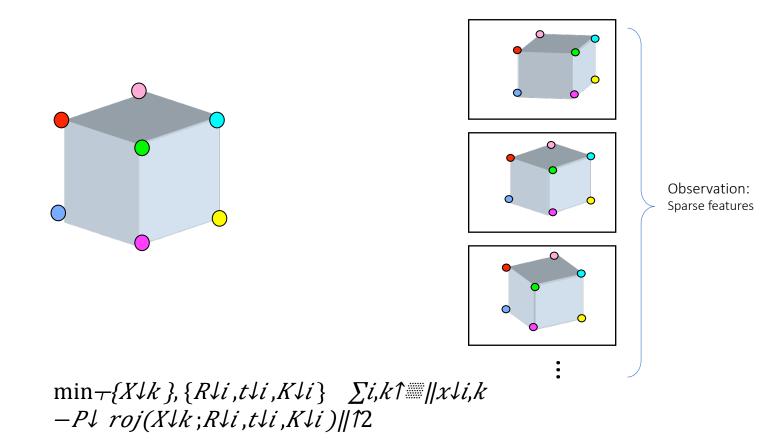
"*Building Rome in a Day*", Sameer Agarwal, Noah Snavely, Ian Simon, Steven M. Seitz and Richard Szeliski. ICCV 2009

- Simultaneous localization and mapping (SLAM)
  - One moving camera
- Google ARCore/Apple ARKit
  mobile camera + IMU
- Bundle Adjustment
  - a technique to simultaneously optimize both geometry and camera poses
- KinectFusion
  - One depth camera, detailed geometry



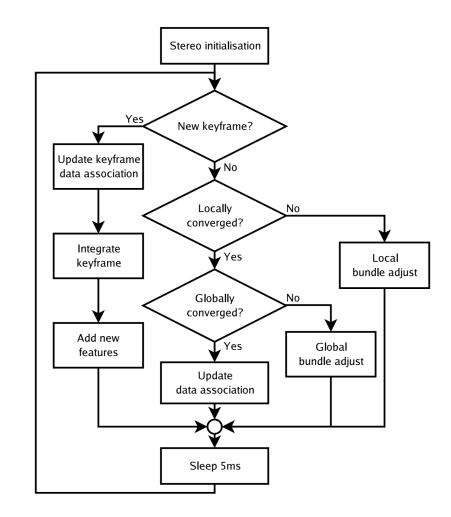
"*KinectFusion: Real-time dense surface mapping and tracking.*" Newcombe, Richard A., Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. ISMAR 2011

# Bundle Adjustment



# SLAM

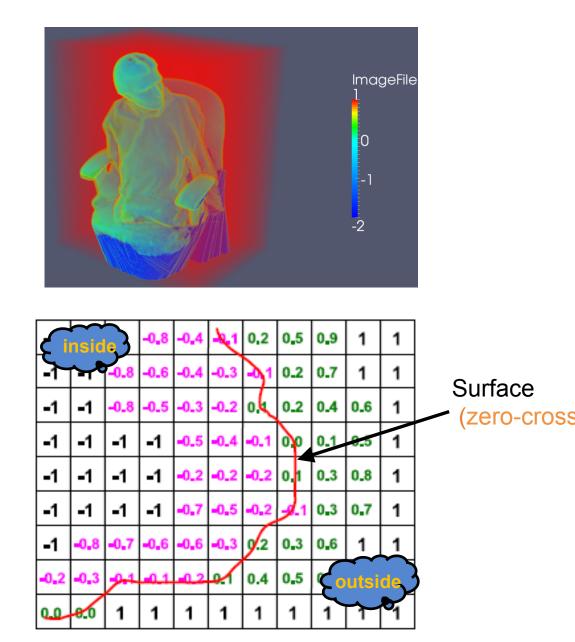
- detect sparse feature points, eg. SIFT
- Initial the system with two-viewgeometry
- Estimate camera poses for later frames by matching 2D features with 3D points



"Parallel Tracking and Mapping for Small AR Workspaces", Georg Klein and David Murray, ISMAR'07

# KinectFusion

- Data accumulation in a TSDF grid
- Camera pose tracking with Iterative Closest Point (ICP)



**TSDF volume grid** 

# People Reconstruction

# The need for temporal consistency





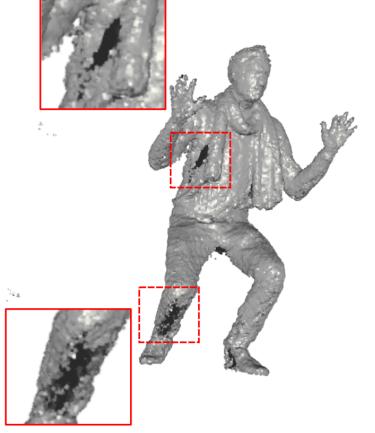
Live Data

**Temporally Consistent Model** 

# The need for temporal consistency



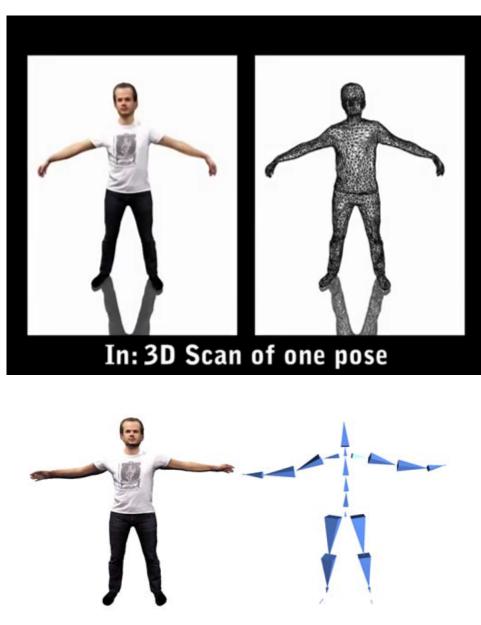
Multiple Point Clouds (Bilaterally-Smoothed)



Fused Live Data (Kinect Fusion)

Temporally consistent model

- Articulated Body tracking
  - Human Body tracking



**Prescan + Skeleton** 

"Marker-less Motion Capture of Skinned Models in a Four Camera Set-up using Optical Flow and Silhouettes", L. Ballan and G. M. Cortelazzo, 3DPVT 2008 33

- Articulated Body tracking
  - Human Body tracking
  - Hand Tracking



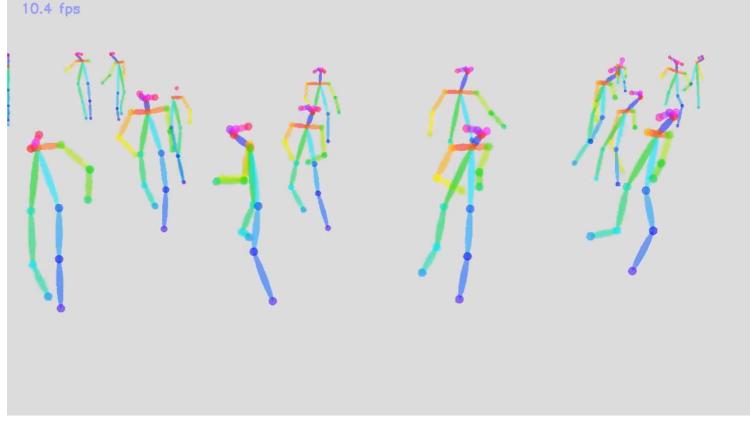
*"Articulated Distance Fields for Ultra-Fast Tracking of Hands Interacting",* Jonathan Taylor, Vladimir Tankovich, Danhang Tang, Cem Keskin, David Kim, Philip Davidson, Adarsh Kowdle, Shahram Izadi. SIGGRAPH Asia 2017

- Articulated Body tracking
  - Human Body tracking
  - Hand Tracking



*"Articulated Distance Fields for Ultra-Fast Tracking of Hands Interacting",* Jonathan Taylor, Vladimir Tankovich, Danhang Tang, Cem Keskin, David Kim, Philip Davidson, Adarsh Kowdle, Shahram Izadi. SIGGRAPH Asia 2017

- Articulated Body tracking
  - Human Body tracking
  - Hand Tracking
- Discriminative method
  - Kinect body tracker
  - deep learning



"*Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields*", Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh, CVPR 2017.

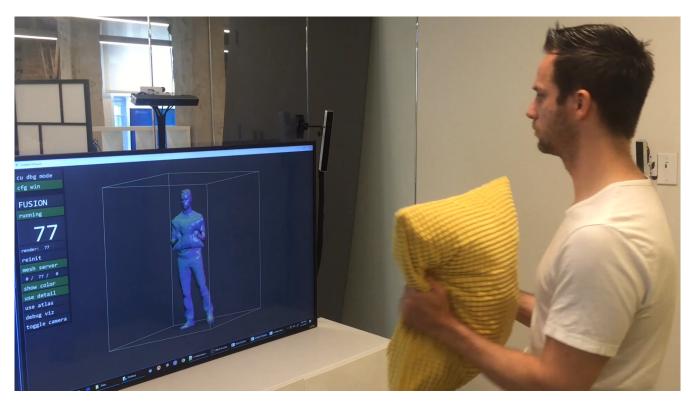
# Reconstruct People

- Articulated Body tracking
  - Human Body tracking
  - Hand Tracking
- Discriminative method
  - Kinect body tracker
  - deep learning
- Face Tracking



Digital Emily Project, Paul Debevec

# General Nonrigid Surface Tracking and Fusion



- Less infra-structure, no green screen.
- No pre-scan, works for general surface

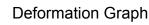


# Non-rigid Alignment

Parameterize motion field as Embedded Deformation Graphs

- One affine transform per ٠ node
- Enforce smoothness over • edges.
- Skinning: attach vertices to neighboring nodes
  - Based on *geodesic* ٠ distance;
  - $\eta_{ine} = \frac{1}{2} k \left[ A(q \downarrow k) (v g \downarrow k) + g \downarrow k + t \downarrow k \right]$





Reference



### mandata $E \downarrow data(G) + \lambda \downarrow reg E \downarrow reg(G)$

geometry alignment learned correspondences texture/color consistency free-space penalization



Sumner, Robert W., Johannes Schmid, and Mark Pauly. "Embedded deformation for shape manipulation." ACM Transactions on Graphics (TOG). Vol. 26. No. 3. ACM, 2007.

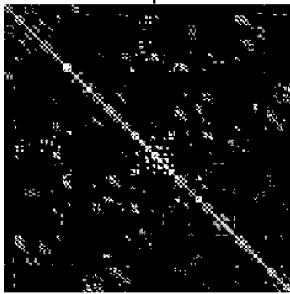
Data

## Solve

## $\mathsf{min}^{\lambda \mathsf{I} \mathsf{data} \, \mathsf{E} \mathsf{I} \mathsf{data} \, (\mathsf{G}) + \lambda \mathsf{I} \mathsf{reg} \, \mathsf{E} \mathsf{I} \mathsf{reg} \, (\mathsf{G})}$

Non-linear least square problem

- Custom/efficient sparse Levenberg-Marquardt solver
- PC GPU ~1ms per iteration
- solve the normal equation at each step:  $(J \uparrow T J + \delta I)h = -J \uparrow T f$



J/T J Non-Zero Pattern

A sparse block matrix, each pixel represents a 7x7 matrix.

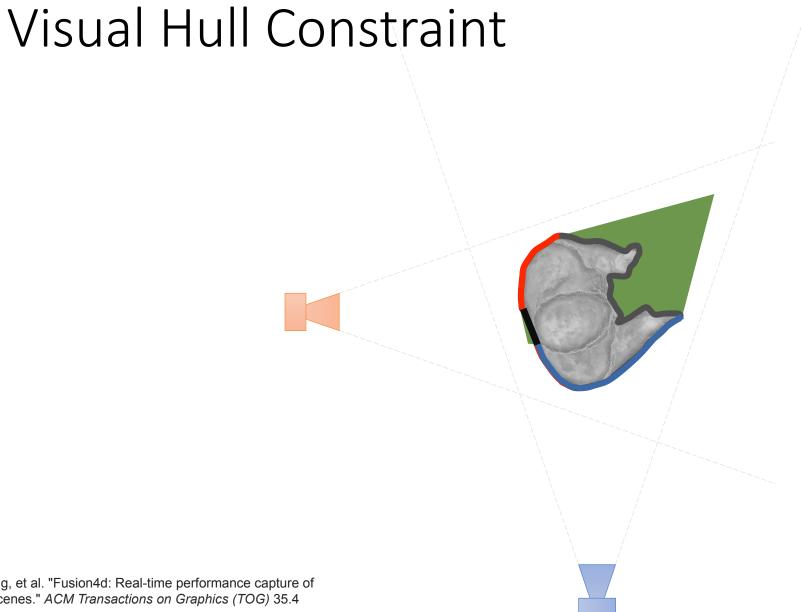
# Data Term $E \downarrow data (G) = \sum m = 1 \uparrow M = \sum n = 1 \uparrow N = (\Psi(v \downarrow m; D \downarrow n)) \uparrow 2$

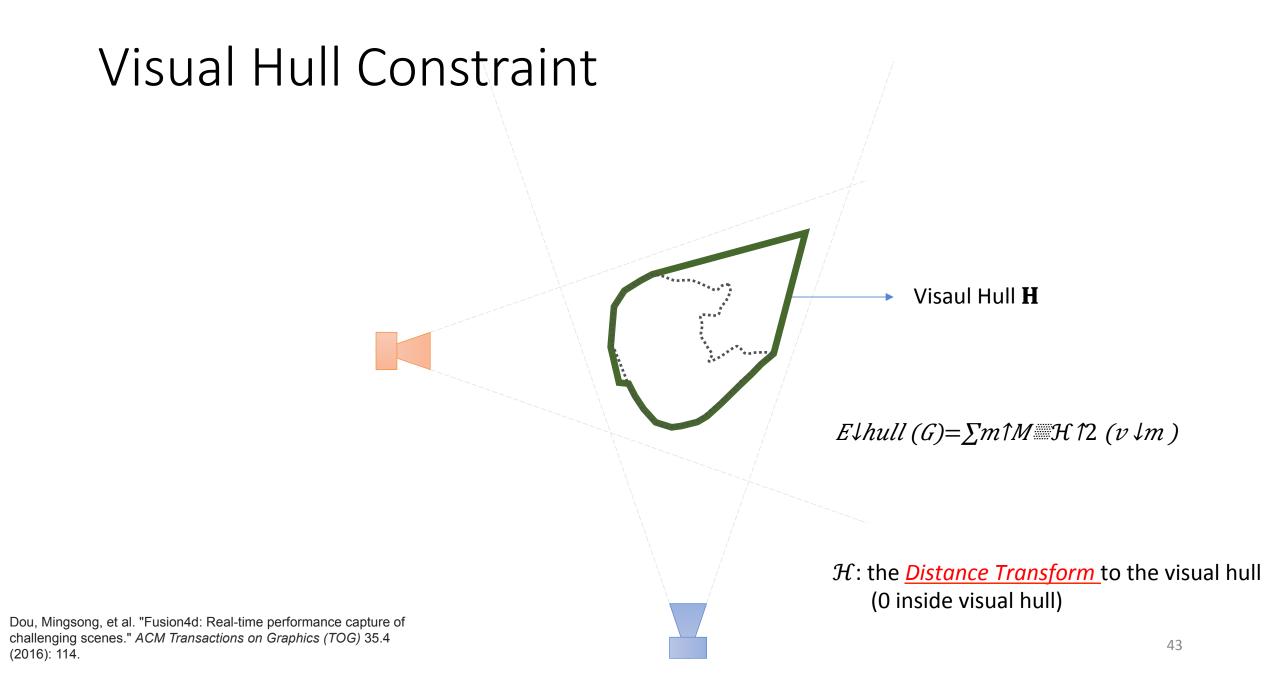
• Alignment residual between model and data

 $\Psi(v \downarrow m; \mathcal{D} \downarrow n) = \delta \downarrow mn n \downarrow m \uparrow T(v \downarrow m - \Gamma \downarrow n (v \downarrow m))$ 

 $\Gamma \downarrow n \ (\nu \downarrow m)$ : projective correspondence of  $\nu \downarrow m$  in depth map  $\mathcal{D} \downarrow n$ 

 $\delta \downarrow mn = \{\blacksquare 1 \text{ if } \nu \downarrow m \text{ is visible to } \mathcal{D} \downarrow n 0 \text{ otherwise} \}$ 





## Color consistency constraint

Measure the difference between the observed color images  $\{I \downarrow k\}$  $\downarrow k=1 \uparrow K$  and the reconstructed image by projecting the per-vertexcolored mesh:





## Single Sensor Results



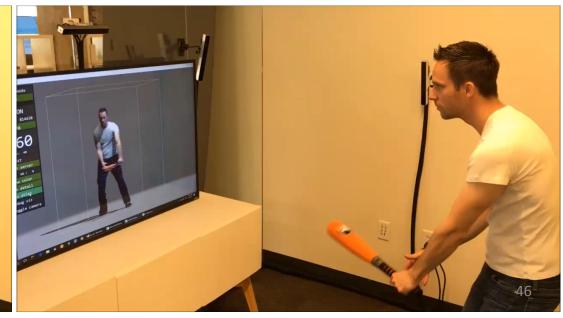




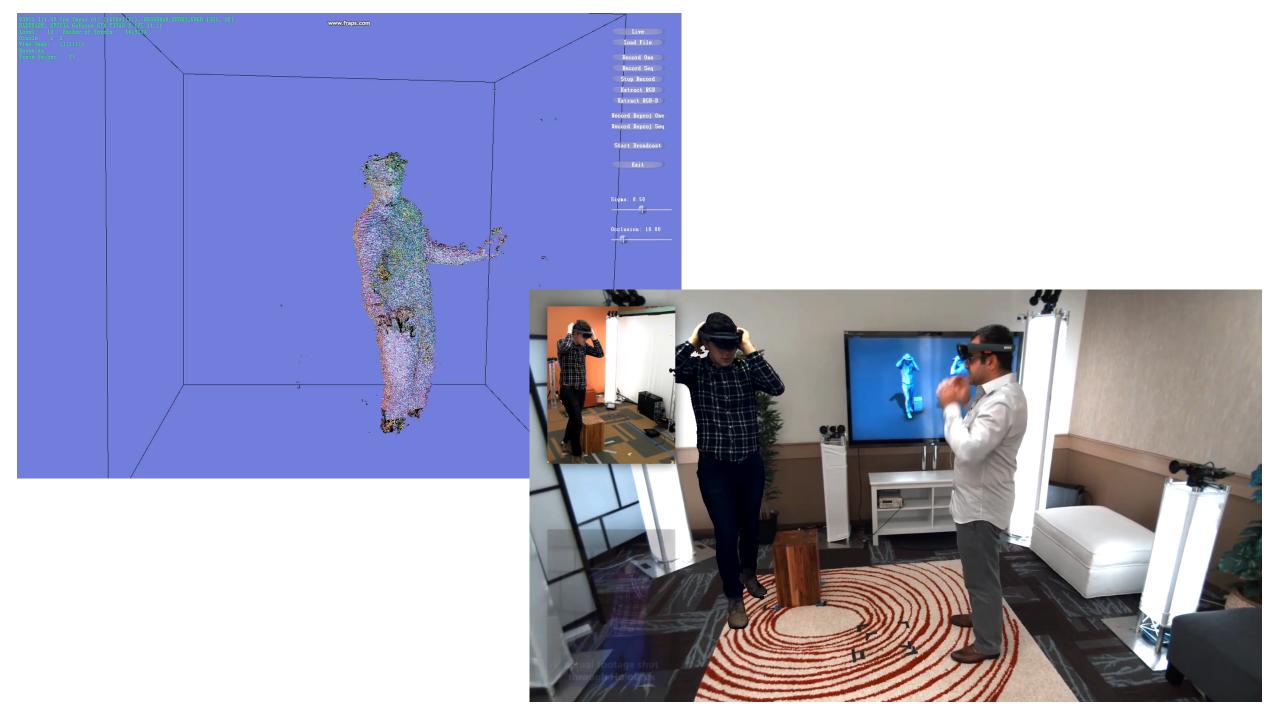
# High Speed 360 Capture





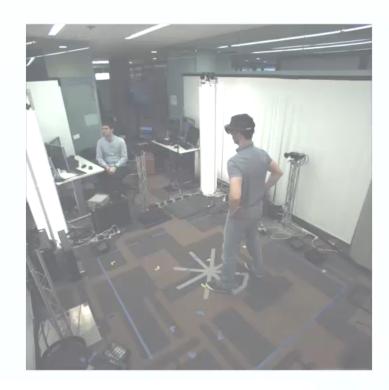


# Normal environment 24 cameras Real-time

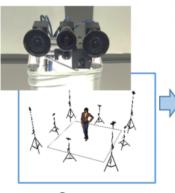








8 Pods

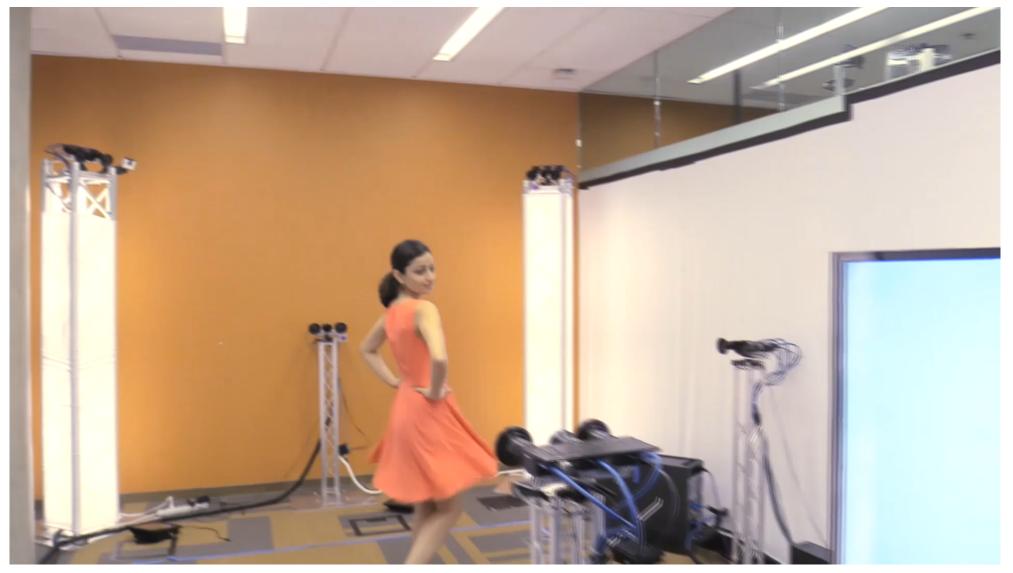


Capture

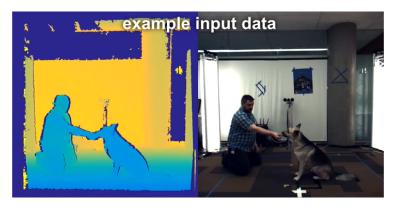
SITE A

SITE B

# Results



# Non-human examples





real-time multi-view reconstruction

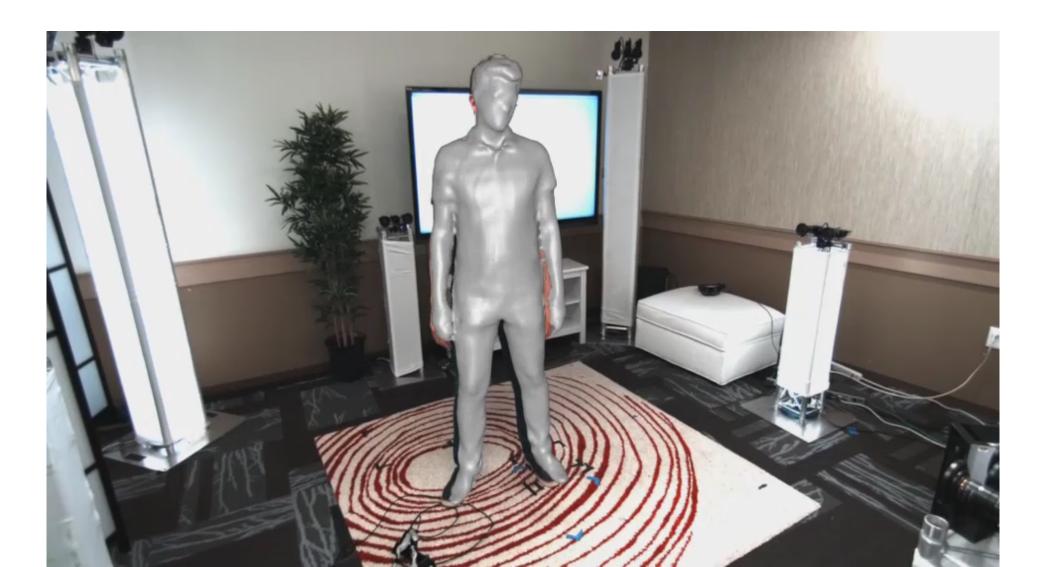
# **Usage Scenarios**

# Challenges for the future: FoV

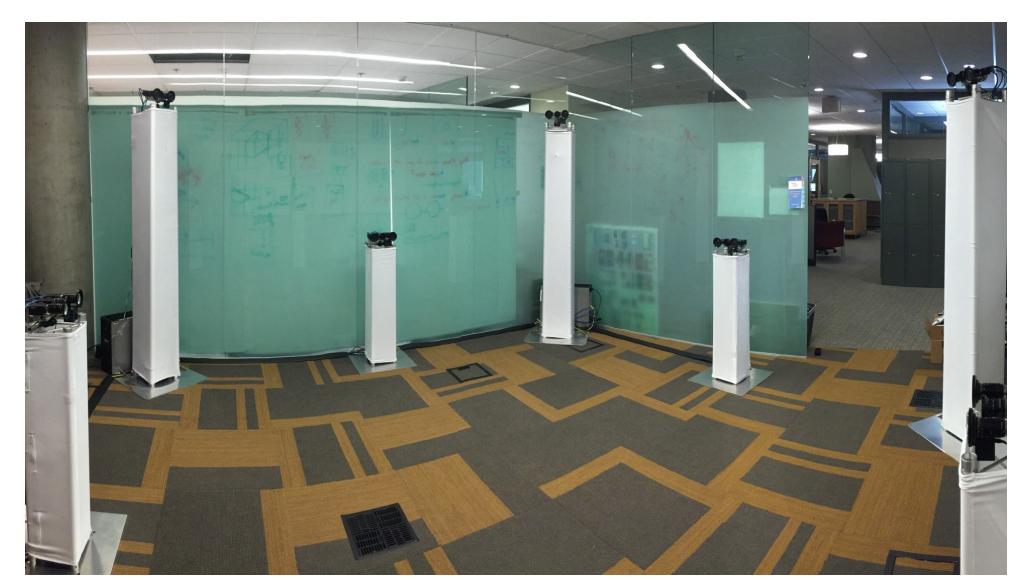




# Challenges: FoV



# Challenges: Infrastructure



# Challenges: Headset removal



## Compression

- Raw Data from Cameras = 23 Gb/s
- Current HoloPort Compression = 1 Gb/s
- HD video = 10 Mb/s
- Exploit temporally consistency and texture atlas

# Thank you!