# GAMES Multi-View Stereo



Qixing Huang August 6<sup>th</sup> 2021



## **Two-View Stereo**

## **Binocular Stereo**

• Given a calibrated binocular stereo pair, fuse it to produce a depth image

image 1





#### Dense depth map



## **Basic Stereo Matching Algorithm**



- For each pixel in the first image
  - Find corresponding epipolar line in the right image
  - Examine all pixels on the epipolar line and pick the best match
  - Triangulate the matches to get depth information
- Simplest case: epipolar lines are corresponding scanlines

   When does this happen?

## Basic stereo matching algorithm



- For each pixel in the first image
  - Find corresponding epipolar line in the right image
  - Examine all pixels on the epipolar line and pick the best match
  - Triangulate the matches to get depth information
- Simplest case: epipolar lines are corresponding scanlines
  - When does this happen?

## Simplest Case: Parallel Images



- Image planes of cameras are parallel to each other and to the baseline
- Camera centers are at same height
- Focal lengths are the same
- Then, epipolar lines fall along the horizontal scan lines of the images



Disparity is inversely proportional to depth!





Conf. Computer Vision and Pattern Recognition, 1999.

## **Rectification Example**



### Correspondence search



- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

### Correspondence search



### Correspondence search



## Effect of window size









W = 20

- -Smaller window
  - + More detail
  - More noise
- -Larger window
  - + Smoother disparity maps
  - Less detail

## Results with window search

Data



#### Window-based matching

Ground truth





## Non-local constraints

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image



## Non-local constraints

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering

- Corresponding points should be in the same order in both views



## Non-local constraints

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering

- Corresponding points should be in the same order in both views



Ordering constraint doesn't hold

### **Consistency Constraints**

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views
- Smoothness
  - We expect disparity values to change slowly (for the most part)

MRF Formulation:

$$E(d) = E_d(d) + \lambda E_s(d)$$
Pixel matching score Consistency Scores

## Comparsion

Window-Based Search:





Ground Truth

Graph Cut:

### Stereo matching as energy minimization



• Graph-cuts can be used to minimize such energy

Y. Boykov, O. Veksler, and R. Zabih, <u>Fast Approximate Energy Minimization via Graph Cuts</u>, PAMI 2001

## **Multi-View Stereo**

## **Multi-View Stereo from Internet Collections**

[Goesele et al. 07]





## Challenges

Appearance variation



Resolution



• Massive collections

82754 results for photos matching notre and dame and paris

## Law of Nearest Neighbors



206 Flickr images taken by 92 photographers









#### 4 best neighboring views











#### reference view





### Local view selection

- Automatically select neighboring views for each point in the image
- Desiderata: good matches AND good baselines









#### 4 best neighboring views











#### reference view





### Local view selection

- Automatically select neighboring views for each point in the image
- Desiderata: good matches AND good baselines









#### 4 best neighboring views











#### reference view



### Local view selection

- Automatically select neighboring views for each point in the image
- Desiderata: good matches AND good baselines

#### Notre Dame de Paris

653 images 313 photographers







129 Flickr images taken by 98 photographers



merged model of Venus de Milo































































































56 Flickr images taken by 8 photographers





#### merged model of Pisa Cathedral



Accuracy compared to laser scanned model: 90% of points within 0.25% of ground truth

## **Use of Structural Priors**
## MVS has been so successful



#### [Apple maps]

High-quality passive facial performance capture using anchor frames [T. Beeler, F. Hahn, D. Bradley, B. Bickel, P. Beardsley, C. Gotsman, M. Gross, 2011]



# What if no texture...



# What if no texture...



# **Existing Method Fails**





#### [Furukawa and Ponce, 2007]

# **Enforce Prior in MVS**

- Manhattan-world assumption
  - Planarity
  - Orthogonality





## Standard Depthmap Markov Random Field (MRF)

$$E = \sum_{p} E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$

A Comparative Study of Energy Minimization Methods for Markov Random Fields with Smoothness-Based Priors [Szeliski et al., PAMI 2008]

 $E = \sum_p E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$ 



$$E = \sum_p E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$



$$E = \sum_p E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$



$$E = \sum_{p} E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$



$$E = \sum_{p} E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$



$$E = \sum_{p} E_d(d_p) + \sum_{\{p,q\} \in \mathbb{N}} E_s(d_p, d_q)$$



 $E = \sum E_d(d_p) + \sum E_s(d_p, d_q)$ p $\{p,q\} \in \mathbb{N}$ 

Graph-cuts (alpha-expansion) gives you very good solutions







### How to enforce piecewise planar

- 1. Advanced MRF and optimization
  - Global Stereo Reconstruction under Second Order Smoothness Priors [Woodford et al., CVPR 2008] Best Paper Award
- 2. Integrate with top-down (primitive) approach
  - Manhattan World Stereo [Furukawa et al., CVPR 2009]
  - Piecewise Planar Stereo for Image-based rendering [Sinha et al., ICCV 2009]
  - Fusion of Feature- and Area-Based Information for Urban Buildings Modeling from Aerial Imagery [Zebedin et al., ECCV 2008]
  - Piecewise Planar and Non-Planar Stereo for urban Scene Reconstruction [Gallup et al., CVPR 2010]

# Advanced MRF for Depthmap Estimation

#### Reference image



#### Ground truth





Global Stereo Reconstruction under Second Order Smoothness Priors [Woodford et al., CVPR 2008] Best Paper Award

# Standard MRF



# Standard MRF

$$E_s(d_p, d_q) = |d_p - d_q|$$



#### Standard MRF

# MRF with a triple clique

$$E_s(d_p, d_q) = |d_p - d_q|$$



$$E_s(d_p, d_q, d_r)$$



# Standard MRF MRF with a triple clique

$$E_s(d_p, d_q) = |d_p - d_q|$$



$$E_s(d_p, d_q, d_r) = |d_p + d_r - 2d_q|$$



# MRF with a triple clique

 $E = \sum E_d(d_p) + \sum E_s(d_p, d_q)$  $\{p,q\} \in \mathbb{N}$ p+  $\sum E_s(d_p, d_q, d_r)$  $\{p,q,r\} \in \mathbb{N}_3$ 

# Optimization becomes a challenge

#### Standard MRF (submodular)



- Graph-cuts (alpha-expansion)
- Belief propagation (TRW)

# MRF with a triple clique (non-submodular)



- QPBO
- Belief propagation (TRW?)

Fusion moves...

# Experimental results

Reference image



Neighboring image



Output depthmap





# Comparative experiment

Reference image



#### Standard MRF



Ground truth



#### MRF with a triple clique



### How to enforce piecewise planar

- 1. Advanced MRF and optimization
  - Global Stereo Reconstruction under Second Order Smoothness Priors [Woodford et al., CVPR 2008] Best Paper Award
- 2. Integrate with top-down (primitive) approach
  - Manhattan World Stereo [Furukawa et al., CVPR 2009]
  - Piecewise Planar Stereo for Image-based rendering [Sinha et al., ICCV 2009]
  - Fusion of Feature- and Area-Based Information for Urban Buildings Modeling from Aerial Imagery [Zebedin et al., ECCV 2008]
  - Piecewise Planar and Non-Planar Stereo for urban Scene Reconstruction [Gallup et al., CVPR 2010]



## Planemap

#### Extract Manhattan directions Extract planes



#### [Furukawa and Ponce, 2007]

## Planemap



## Planemap



Depthmap to Planemap Depthmap  $(d_p)$  $E = \sum E_d(d_p) + \sum E_s(d_p, d_q)$  $\{p,q\} \in \mathbb{N}$ Planemap  $(h_p)$  $E = \sum E_d(h_p) + \sum E_s(h_p, h_q)$  $\{p,q\} \in \mathbb{N}$ p



Planemap MRF  $E = \sum E_d(h_p) + \sum E_s(h_p, h_q)$  $\{p,q\} \in \mathbb{N}$ p



Planemap MRF  $E = \sum E_d(h_p) + \sum E_s(h_p, h_q)$  $\{p,q\} \in \mathbb{N}$ p



Planemap MRF  $E = \sum E_d(h_p) + \sum E_s(h_p, h_q)$  $\{p,q\} \in \mathbb{N}$ p


Planemap MRF  $E = \sum E_d(h_p) + \sum E_s(h_p, h_q)$  $\{p,q\} \in \mathbb{N}$ p



# Comparison



#### Standard method

Manhattan Planemap

## **Reconstruction Results**

Kitchen - 22 images

house - 148 images

gallery - 492 images











## How to enforce piecewise planar

- 1. Advanced MRF and optimization
  - Global Stereo Reconstruction under Second Order Smoothness Priors [Woodford et al., CVPR 2008] Best Paper Award
- 2. Integrate with top-down (primitive) approach
  - Manhattan World Stereo [Furukawa et al., CVPR 2009]
  - Piecewise Planar Stereo for Image-based rendering [Sinha et al., ICCV 2009]
  - Fusion of Feature- and Area-Based Information for Urban Buildings Modeling from Aerial Imagery [Zebedin et al., ECCV 2008]
  - Piecewise Planar and Non-Planar Stereo for urban Scene Reconstruction [Gallup et al., CVPR 2010]

# **Relaxing Mahnattan**



- Use sparse lines + sparse points to detect planes
- MRF + Graph-cuts

Piecewise Planar Stereo for Image-based rendering [Sinha et al., ICCV 2009]

## **Relaxing Manhattan**



[Sinha et al.]





# Relaxing Manhattan



[Sinha et al.]





# Enable Curved Surfaces

• Building reconstruction from a top down view



Fusion of Feature- and Area-Based Information for Urban Buildings Modeling from Aerial Imagery [Zebedin et al., ECCV 2008]

# Enable Curved Surfaces





## How to enforce piecewise planar

- 1. Advanced MRF and optimization
  - Global Stereo Reconstruction under Second Order Smoothness Priors [Woodford et al., CVPR 2008] Best Paper Award
- 2. Integrate with top-down (primitive) approach
  - Manhattan World Stereo [Furukawa et al., CVPR 2009]
  - Piecewise Planar Stereo for Image-based rendering [Sinha et al., ICCV 2009]
  - Fusion of Feature- and Area-Based Information for Urban Buildings Modeling from Aerial Imagery [Zebedin et al., ECCV 2008]
  - Piecewise Planar and Non-Planar Stereo for urban Scene Reconstruction [Gallup et al., CVPR 2010]

## **3D Reconstruction from Video**

**Street-Side Video** 

# FDC

#### **Real-Time Stereo**









#### Enforce planarity where it looks like a building



### First step: Planar and Non-Planar Stereo

Video



#### Labels = { $\pi_1, \cdots, \pi_N, \pi_\infty, non-plane, discard$ }



Real-Time Stereo

**Plane Detection** 

Planemap (w/ non-planar))



## Second step: Appearance Prior

#### **Video Frames**



## Algorithm

Video



#### Planar/Non-Planar Classification



planar

non-planar

Planar and Non-Planar MRF





**Real-Time Stereo** 



**Plane Detection** 



#### Large-Scale Multi-View Stereo

## Input images



[ diydrones.com ]

## Input images



[ diydrones.com ]

## **Divide and Reconstruct**



#### [ diydrones.com ]

# Large-Scale MVS for Unorganized Photos (Internet Photos)

"Building Rome in a Day" [ Agarwal, Snavely, et al., 2009 ]







## Image-clustering based on a model

"Towards Internet-scale Multi-View Stereo" [furukawa et al., 2010]



## One way to look at the clustering problem



## Formulation



SFM points  $\{P_1, P_2, \ldots\}$ 

Images  $\{I_1, I_2, ...\}$ 

Image clusters  $\{C_1, C_2, ...\}$ 

*P1* is reconstructed well in *C1* based on (image resolutions/distribution)

 $\rightarrow$  "P1 is covered by C1."

"P4 is covered by C2."

## Formulation



SFM points  $\{P_1, P_2, \ldots\}$ 

Images  $\{I_1, I_2, ...\}$ 

Image clusters  $\{C_1, C_2, ...\}$ 

We want to

- 1. remove redundant images
- 2. keep each cluster small (memory limit)
- 3. "cover" many points

## Formulation





Minimize
$$\sum_{k} |C_k| \text{ subject to}$$
(compactness)•  $\forall k \ |C_k| \le \alpha,$ (size)•  $\forall i \ \frac{\{\text{\# of covered points in } I_i\}}{\{\text{\# of points in } I_i\}} \ge \delta.$ (coverage)

1. Start with all the images in a cluster



Minimize
$$\sum_{k} |C_k|$$
 subject to(compactness)•  $\forall k \models C_k \models \leq \alpha$ ,(size)•  $\forall i \quad \frac{\{\text{\# of covered points in } I_i\}}{\{\text{\# of points in } I_i\}} \geq \delta$ .(coverage)

- 1. Start with all the images in a cluster
- 2. Optimize *compactness* while keeping *coverage*





- 1. Start with all the images in a cluster
- 2. Optimize *compactness* while keeping *coverage*
- 3. If *size* constraint is broken, split a cluster





- 1. Start with all the images in a cluster
- 2. Remove redundant images while keeping *coverage*
- 3. If *size* constraint is broken, split a cluster



4. Add an image to each cluster to satisfy *coverage* 



- 1. Start with all the images in a cluster
- 2. Remove redundant images while keeping *coverage*
- 3. If *size* constraint is broken, split a cluster



- 4. Add an image to each cluster to satisfy *coverage*
- 5. Repeat (3,4) until *size* and *coverage* are satisfied





## St. Peter's Basilica

- 1275 images
- 4 clusters
- 6M 3D points




