

# Depth Extraction from Video Using Non- parametric Sampling

Kevin Karsch

University of Illinois

Ce Liu

Microsoft Research  
New England

Sing Bing Kang

Microsoft Research

# Problem Statement

Given an image/video, estimate distance from the camera

No parallax necessary

Camera motion OK

Scene motion OK

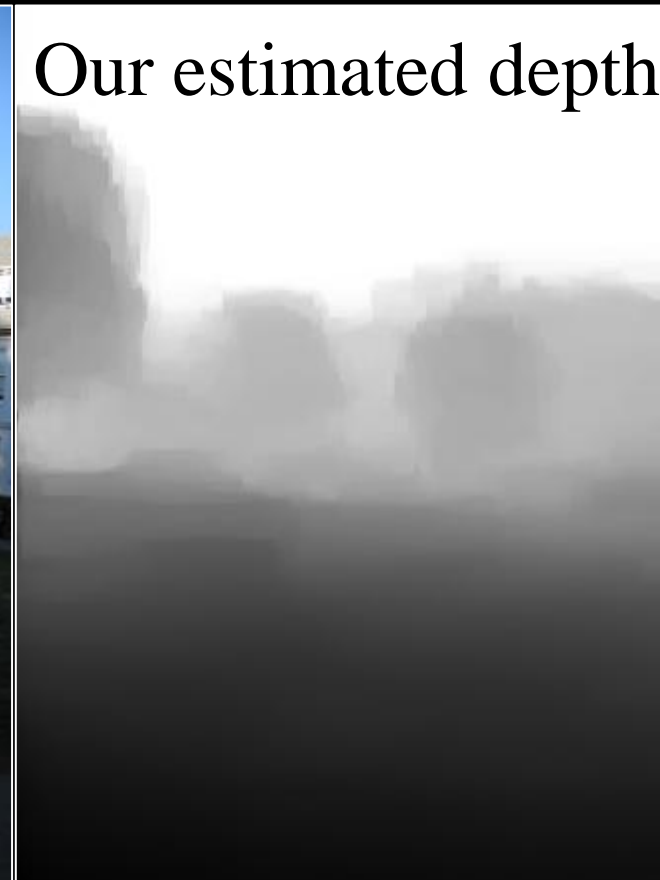
# Problem Statement

Given an image/video, estimate distance from the camera

No parallax necessary

Camera motion OK

Scene motion OK



# Problem Statement

Given an image/video, estimate distance from the camera

No parallax necessary

Camera motion OK

Scene motion OK

Input

Estimated depth



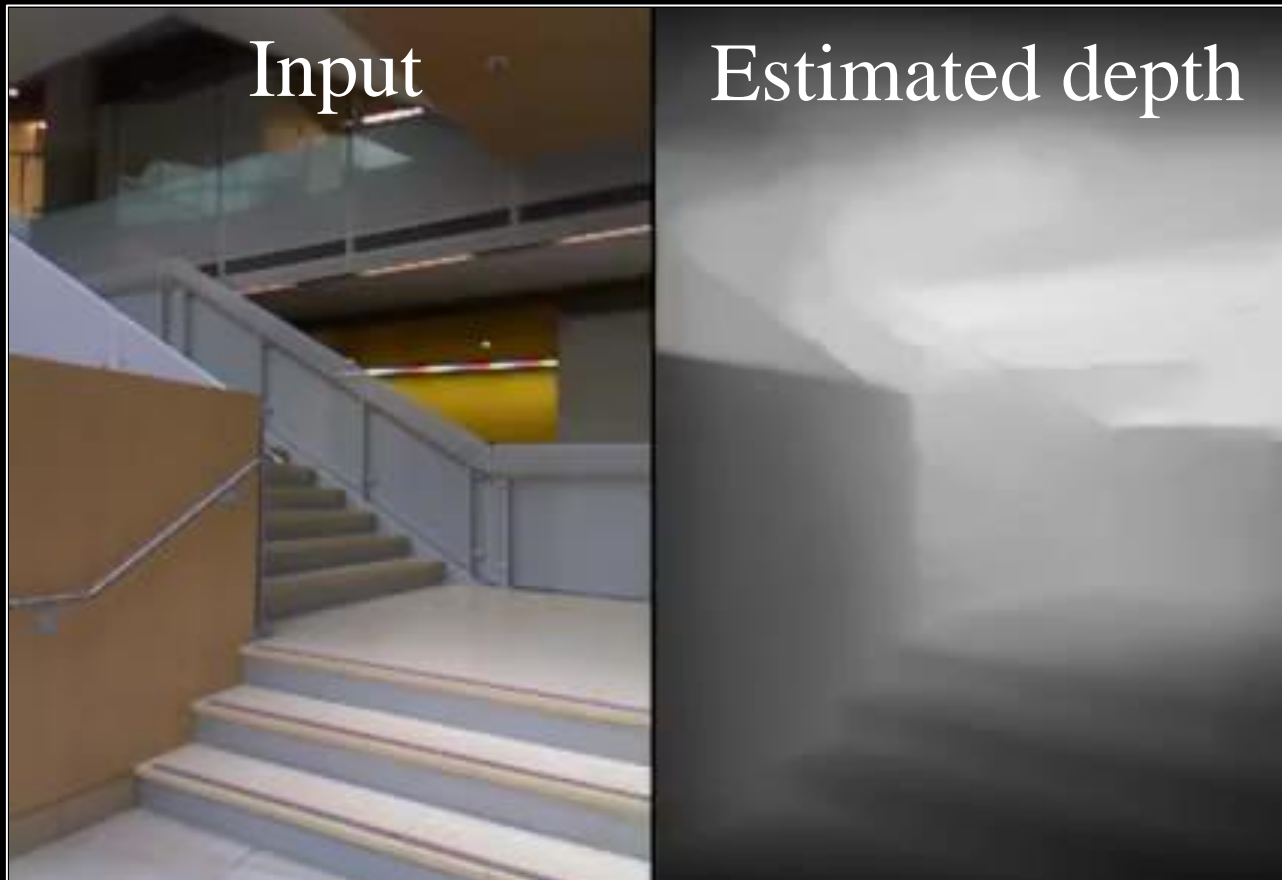
# Problem Statement

Given an image/video, estimate distance from the camera

No parallax necessary

Camera motion OK

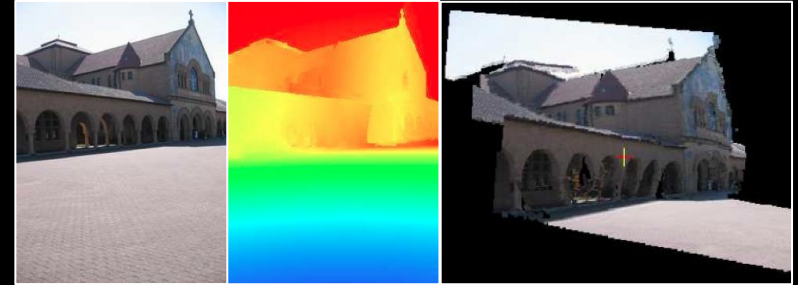
Scene motion OK



# Related Work



[Zhang et al. '09]



[Liu et al. '10]

## Multiview reconstruction

- Very accurate for videos with moving camera
- May fail for dynamic scenes

[Newcombe and Davidson '10]

[Furukawa and Ponce '09]

[Zhang et al. '09]

...

## Parametric learning

- Works well for single images
- No literature on extending to video

[Liu et al. '10]

[Saxena et al. '09]

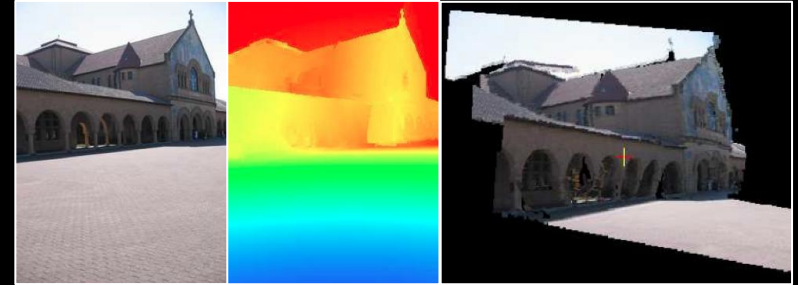
[Hoiem et al. '05]

...

# Related Work



[Zhang et al. '09]



[Liu et al. '10]

## Multiview reconstruction

- Very accurate for videos with moving camera
- May fail for dynamic scenes

[Newcombe and Davidson '10]

[Furukawa and Ponce '09]

[Zhang et al. '09]

...

## Parametric learning

- Works well for single images
- No literature on extending to video

[Liu et al. '10]

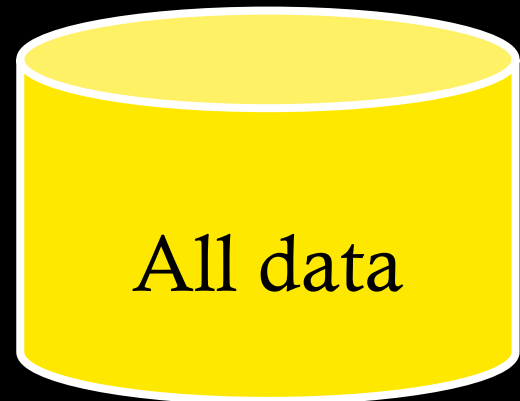
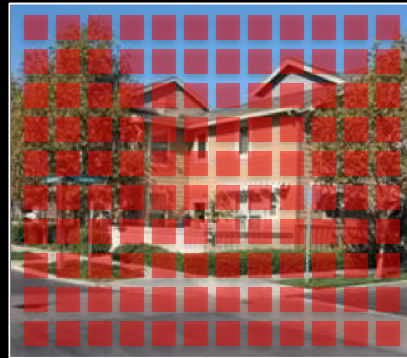
[Saxena et al. '09]

[Hoiem et al. '05]

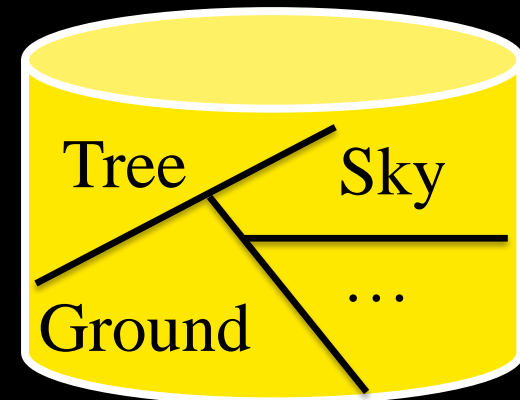
...

# Training set

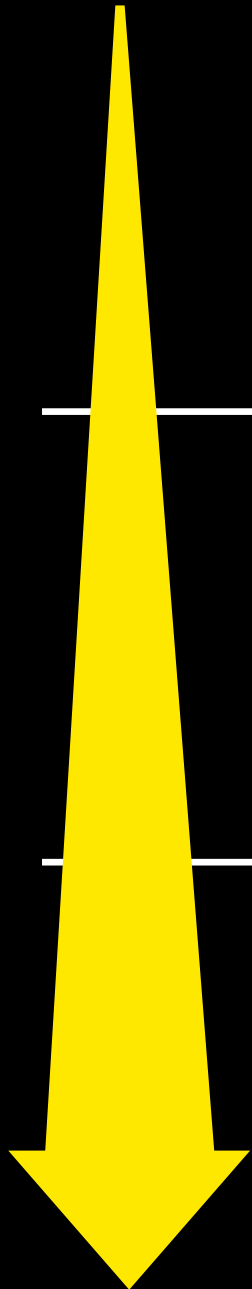
Pixel level [Saxena et al. '05]



Object level [Lui et al. '10]

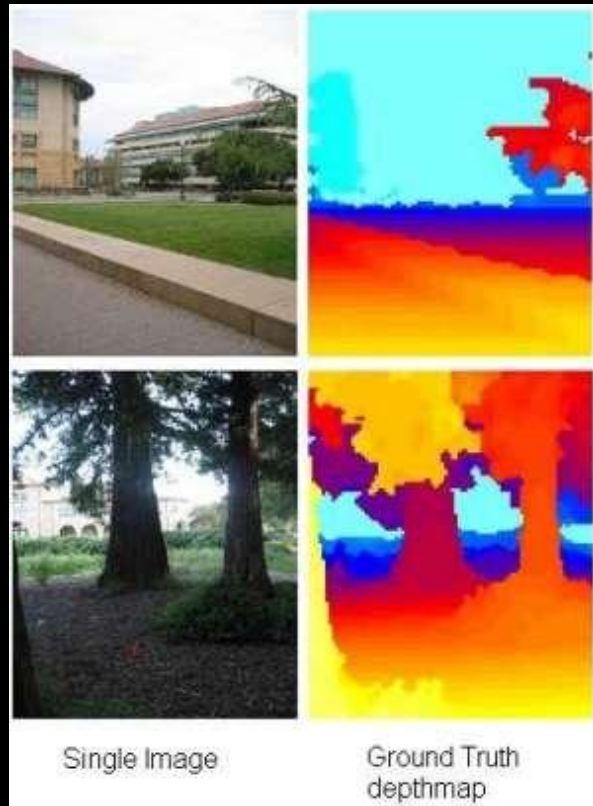


Scene level (ours)





# RGBD Datasets



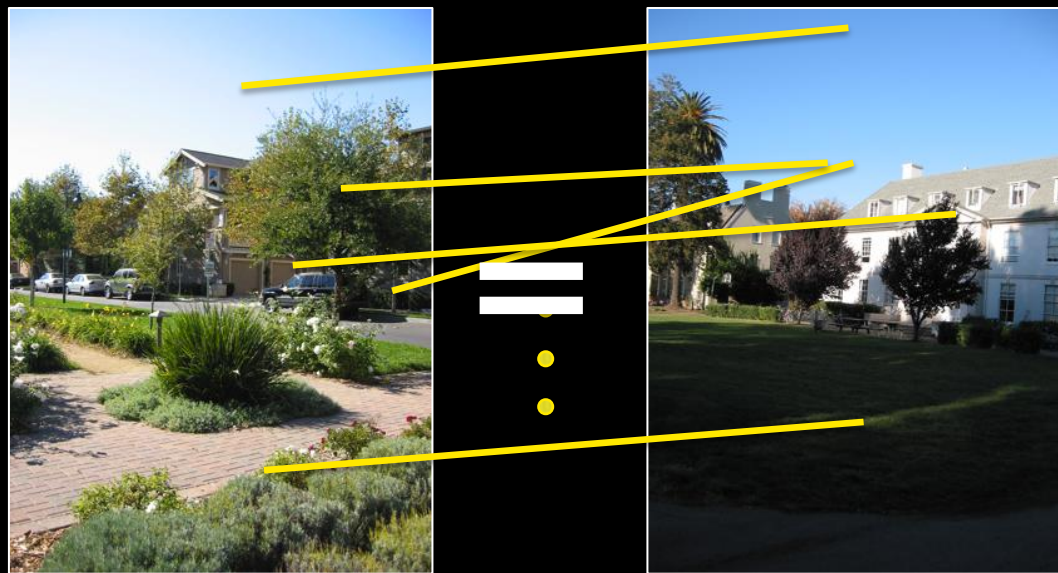
Laser rangefinder  
*Outdoor scenes*  
[Saxena et al.]



MSR-V3D  
*Indoor scenes*  
(Ours)

# SIFT Flow Refresher

- Optical flow using dense SIFT features
  - Larger search window
  - Modified smoothness constraints
- Scenes rearranged so semantics are matched



A

B

# SIFT Flow Refresher

- Optical flow using dense SIFT features
  - Larger search window
  - Modified smoothness constraints
- Scenes rearranged so semantics are matched

$\Psi$   
Warping  
operator



A



B

# SIFT Flow Refresher

- Optical flow using dense SIFT features
  - Larger search window
  - Modified smoothness constraints
- Scenes rearranged so semantics are matched

$\Psi$   
Warping  
operator



A

=



B

# SIFT Flow Refresher

- Optical flow using dense SIFT features
  - Larger search window
  - Modified smoothness constraints
- Scenes rearranged so semantics are matched

$\Psi$   
Warping  
operator



A

=



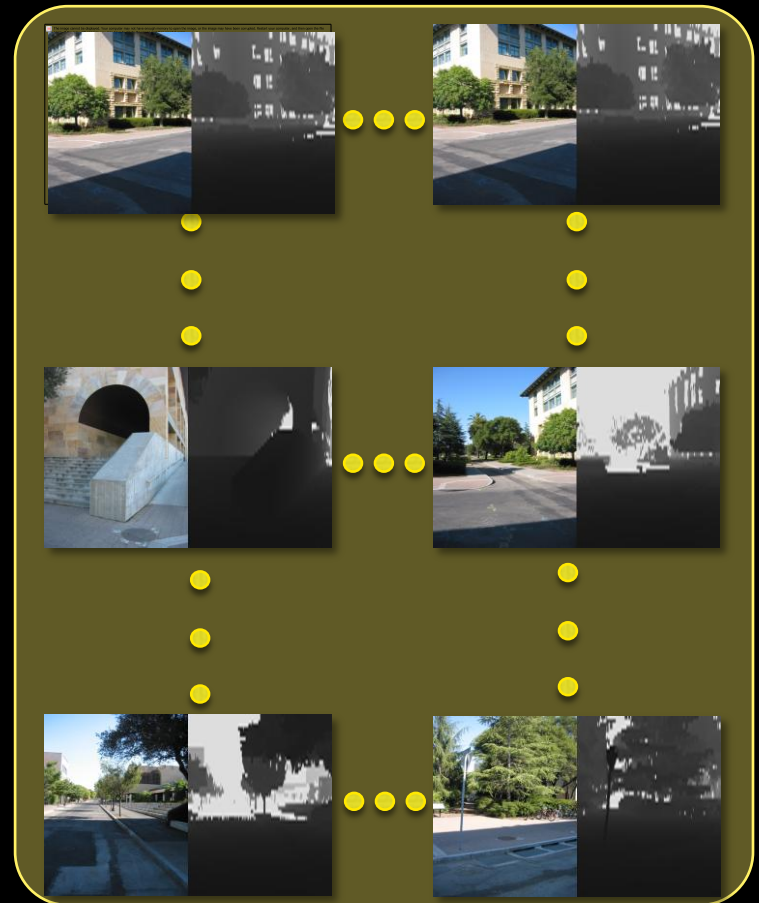
B

# Algorithm

Input image



RGBD Database

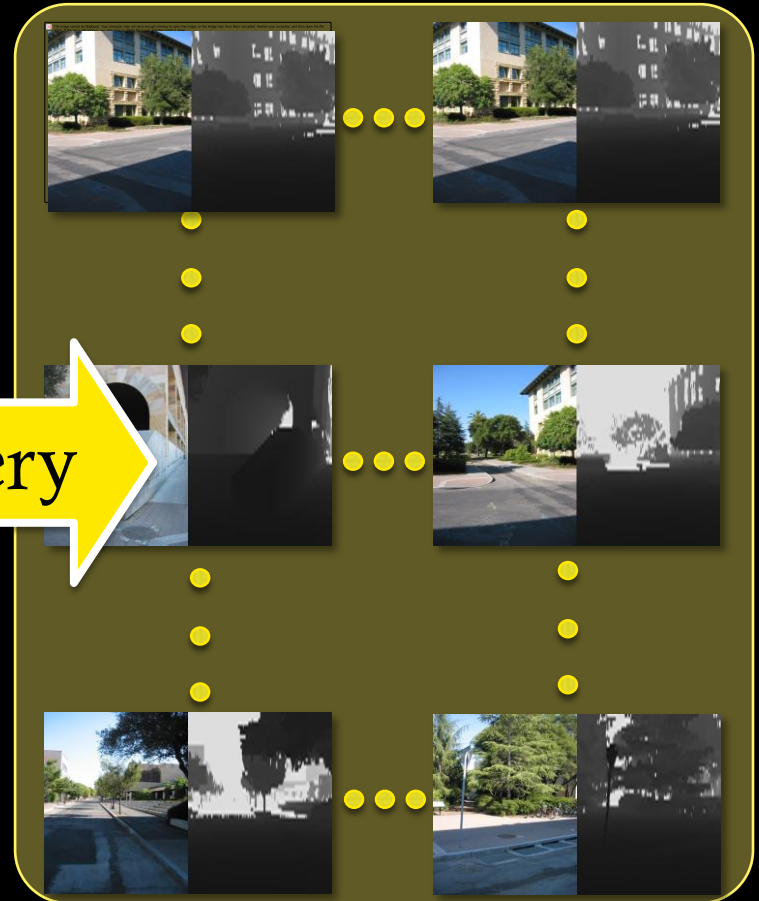


# Algorithm

Input image



RGBD Database

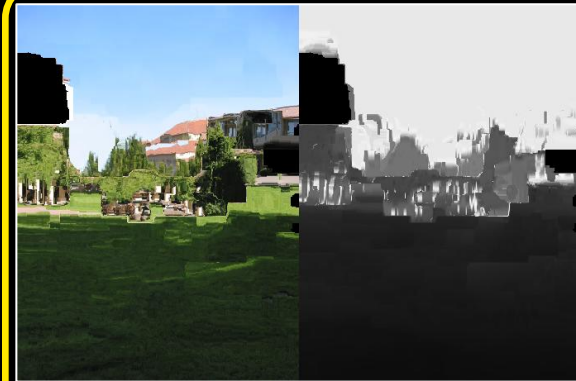
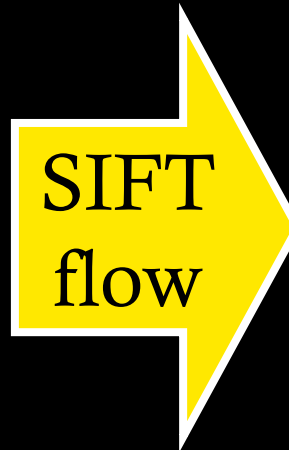
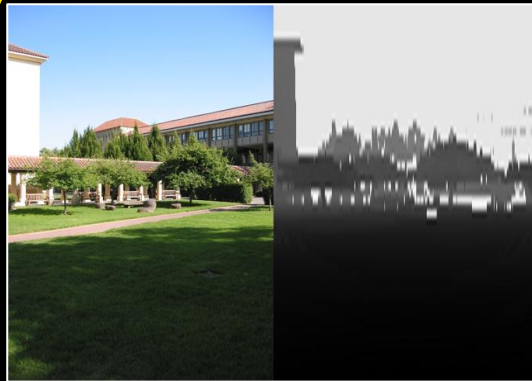


# Algorithm

Candidates

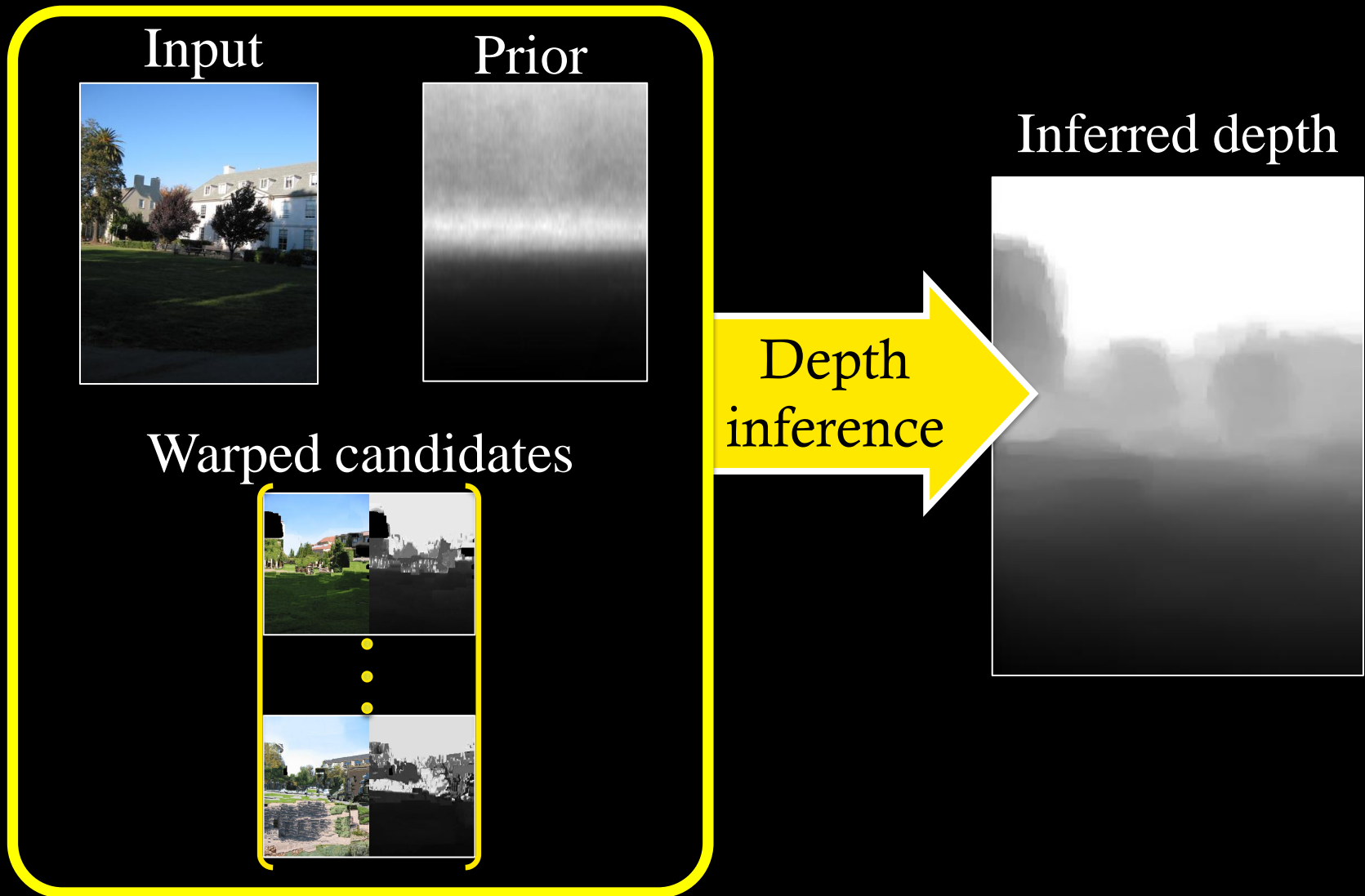
Warped candidates

Input





# Algorithm



# Inference

$$\operatorname{argmin}_D E(D) = \text{Enforce depth to match candidates}$$
$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right]$$

$$+ \underbrace{\alpha s_i |\nabla D_i|_1}_{\text{Spatial smoothness}} + \underbrace{\beta |D_i - \text{prior}_i|_1}_{\text{Match to database mean}}$$

$D$  : inferred depth

$C$  : warped  
candidate depth

$w$  : depth confidence

$S$  : image-based  
weights

$\alpha, \beta, \gamma$  : constant  
weights

- Both *absolute* and *relative* depth are transferred
- Regularize with smoothness and prior

# Inference

argmin <sub>$D$</sub>   $E(D)$  = Enforce depth to match candidates

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right]$$

Not a discrete MRF!

$D$  : inferred depth  
 $C$  : warped candidate depth  
 $w$  : depth confidence  
 $S$  : image-based weights  
 $\alpha, \beta, \gamma$  : constant weights

Spatial smoothness

Match to database mean

- Both *absolute* and *relative* depth are transferred
- Regularize with smoothness and prior

Input

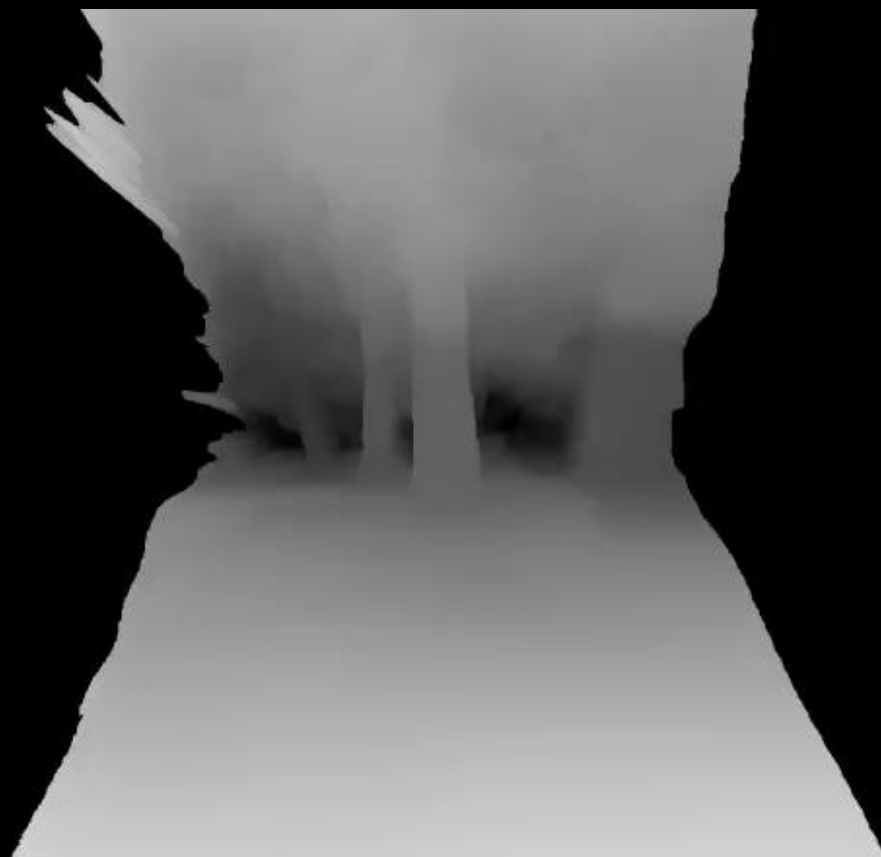


True depth



Inferred depth





Input

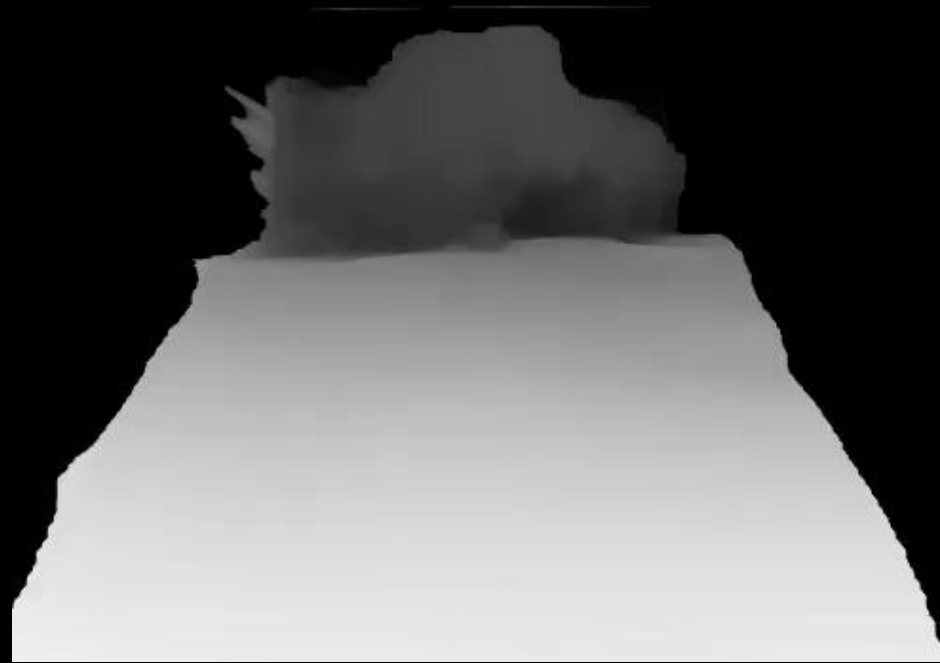


True depth



Inferred depth





$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right]$$

$$+ \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Result *without* relative  
depth term ( $\gamma = 0$ )



Result *with* relative  
depth term ( $\gamma > 0$ )





$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right]$$

$$+ \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Result *without* relative  
depth term ( $\gamma = 0$ )

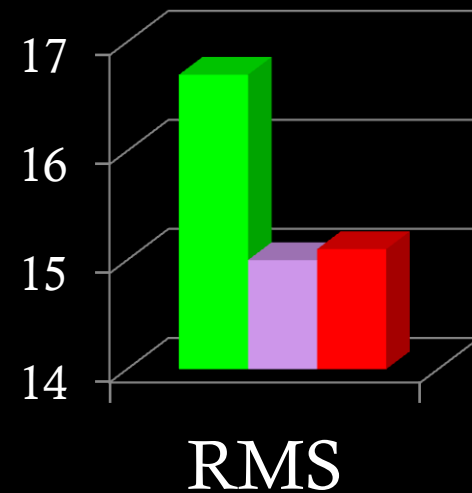
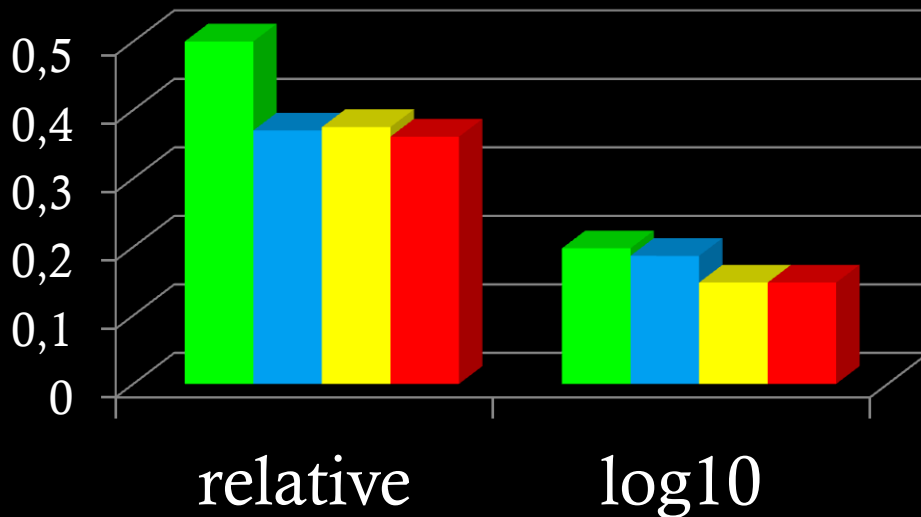


Result *with* relative  
depth term ( $\gamma > 0$ )



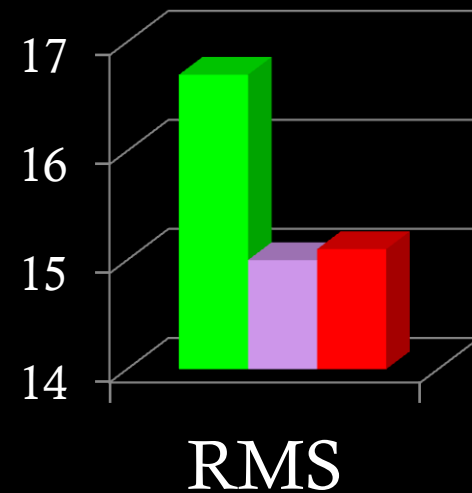
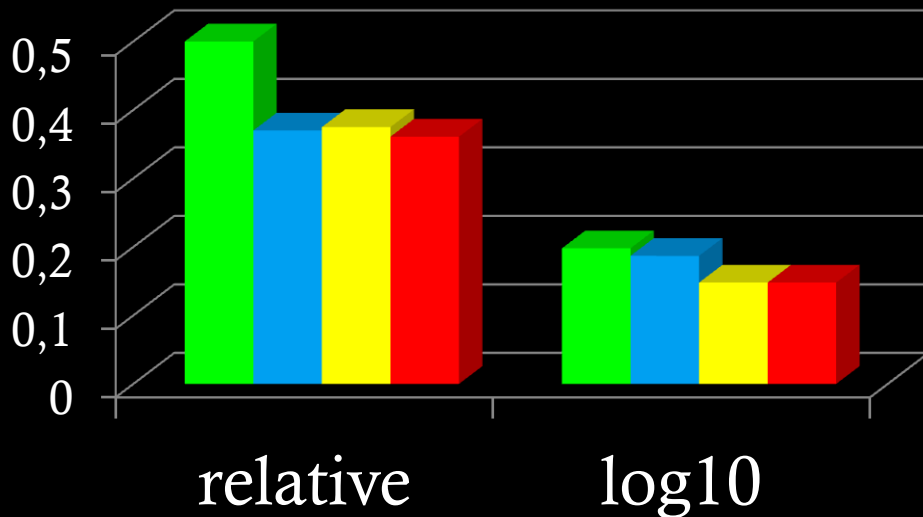
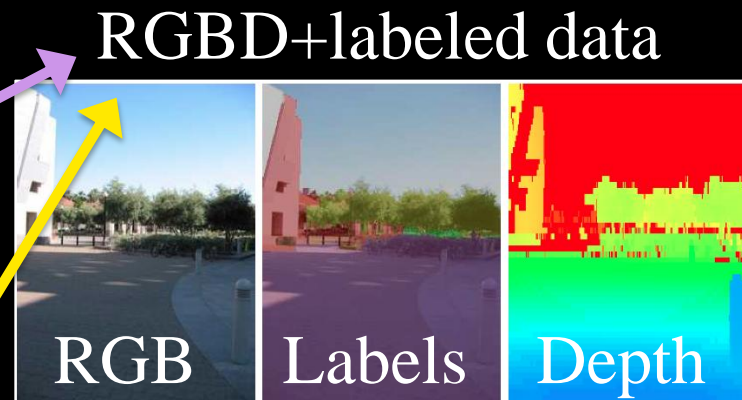
# Evaluation: Make3D Dataset

Method	
	Depth MRF [Saxena et al. '05]
	Make3D [Saxena et al. '09]
	$\theta$ -MRF [Li et al. '11]
	Semantic Labels [Liu et al. '10]
	Depth Transfer (ours)

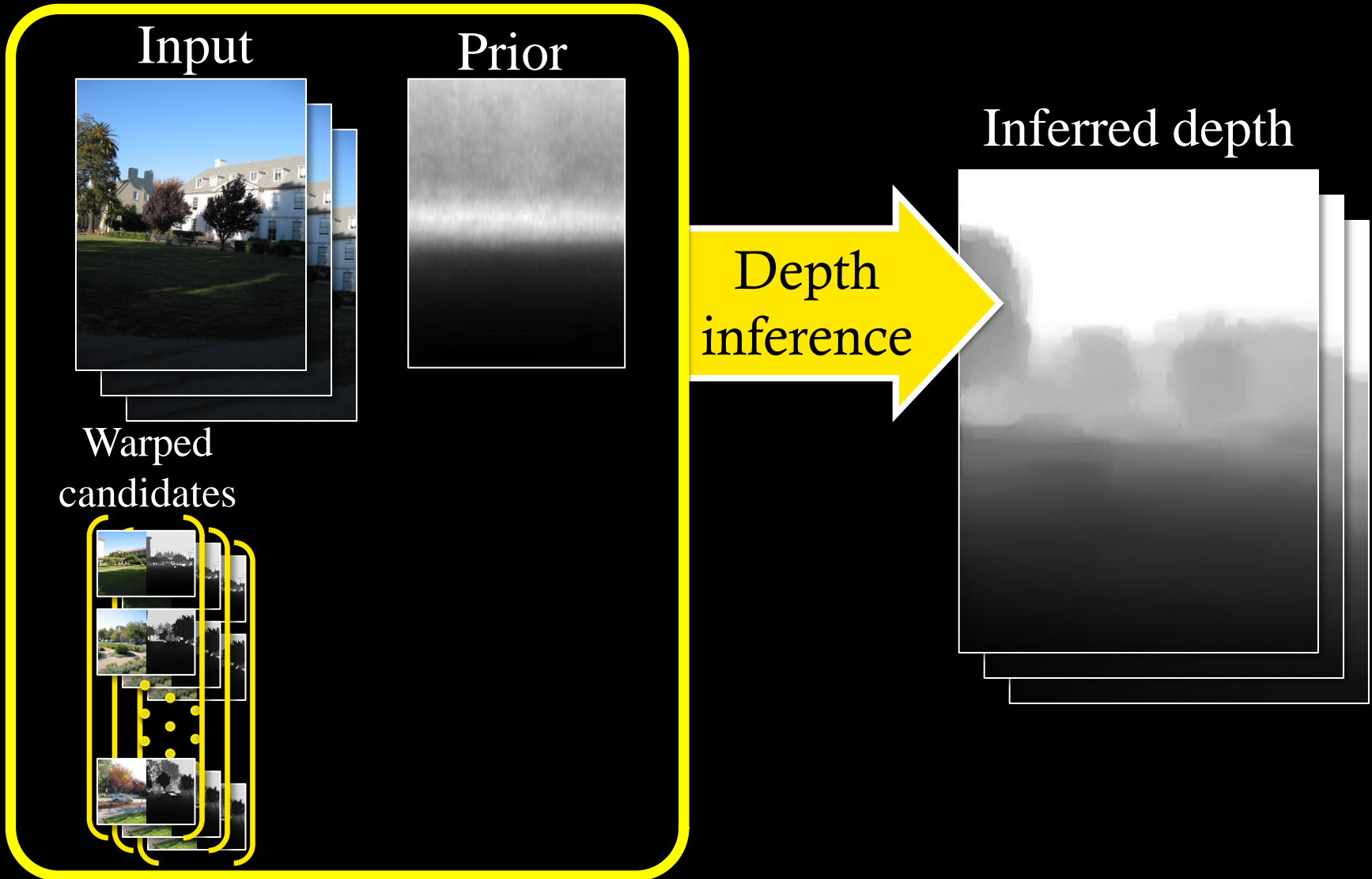


# Evaluation: Make3D Dataset

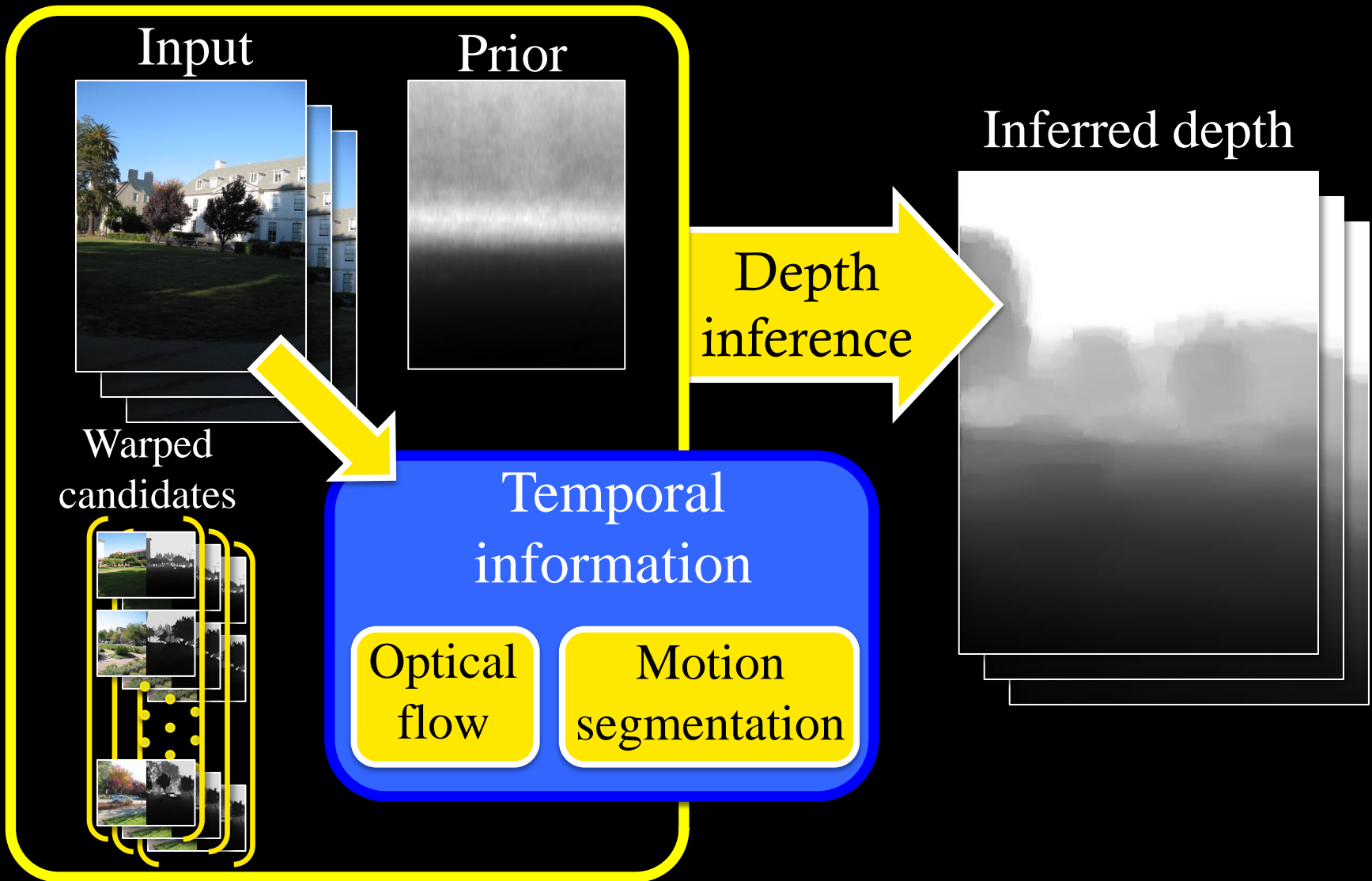
Method
Depth MRF [Saxena et al. '05]
Make3D [Saxena et al. '09]
$\theta$ -MRF [Li et al. '11]
Semantic Labels [Liu et al. '10]
Depth Transfer (ours)



# Video Extension



# Video Extension



# Video Inference

$$\operatorname{argmin}_D E_{\text{video}}(D) =$$

$$\underbrace{E(D)}_{\text{Single image objective}} + \sum_{i \in \text{pixels}} \zeta t_i |\nabla_{\text{flow}} D_i|_1 + \eta m_i |D_i - \mathcal{M}_i|_1$$

$m$  : binary motion mask  
 $\mathcal{M}$  : hypothesized depth of motion mask  
 $\zeta, \eta$  : constant weights

- Depth changes are gradual frame-to-frame
- Moving objects are usually on the ground

# Video Inference

$$\operatorname{argmin}_D E_{\text{video}}(D) =$$

$$\underbrace{E(D)}_{\text{Single image objective}} + \sum_{i \in \text{pixels}} \underbrace{\zeta t_i |\nabla_{\text{flow}} D_i|_1 + \eta m_i |D_i - \mathcal{M}_i|_1}_{\text{Smooth along direction of optical flow}}$$

$m$  : binary motion mask  
 $\mathcal{M}$  : hypothesized depth of motion mask  
 $\zeta, \eta$  : constant weights

- **Depth changes are gradual frame-to-frame**
- Moving objects are usually on the ground

# Video Inference

$$\operatorname{argmin}_D E_{\text{video}}(D) =$$

$$\underbrace{E(D)}_{\text{Single image objective}} + \sum_{i \in \text{pixels}} \underbrace{\zeta t_i |\nabla_{\text{flow}} D_i|_1}_{\text{Smooth along direction of optical flow}} + \underbrace{\eta m_i |D_i - \mathcal{M}_i|_1}_{\text{Coerce moving objects to be "grounded"}}$$

$m$  : binary motion mask  
 $\mathcal{M}$  : hypothesized depth of motion mask  
 $\zeta, \eta$  : constant weights

- Depth changes are gradual frame-to-frame
- **Moving objects are usually on the ground**
  - Motion mask = threshold flow-weighted, relative pixel differences
  - Ce Liu's optical flow <http://people.csail.mit.edu/celiu/OpticalFlow>



Input

Inferred depth



without  
temporal info

with  
temporal info



# Results







# MSR-V3D evaluation

Input

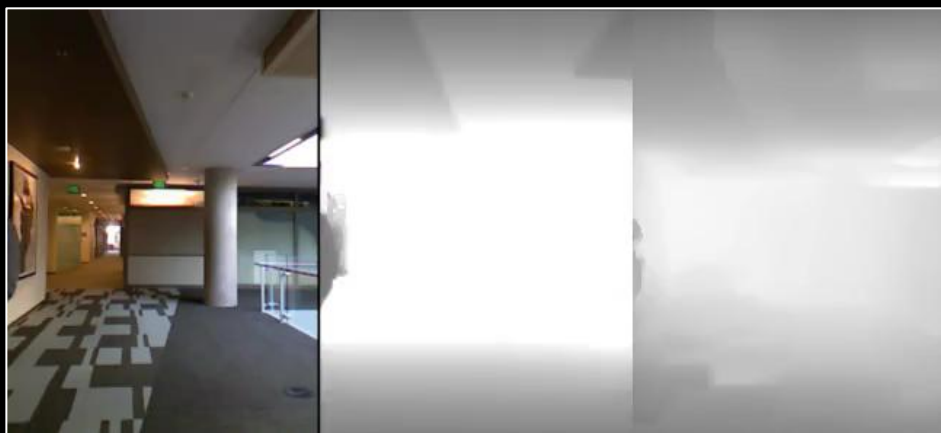
Kinect\*

Ours

Input

Kinect\*

Ours



\*Naïve hole filling applied to Kinect data (for visualization only)

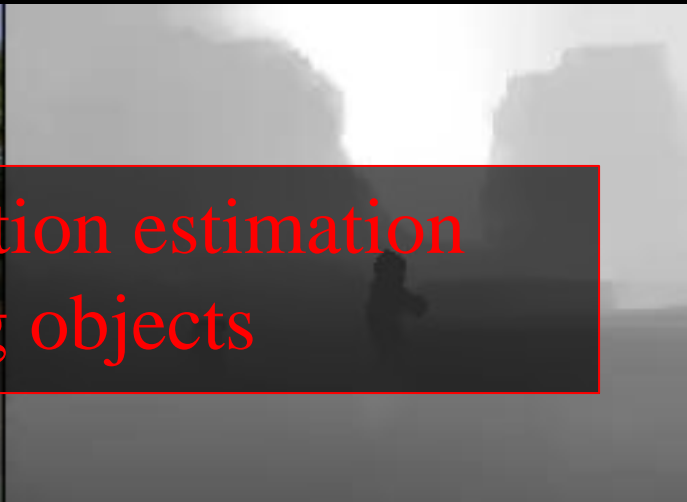
# Limitations



No similar training images



Inaccurate motion estimation  
Floating objects

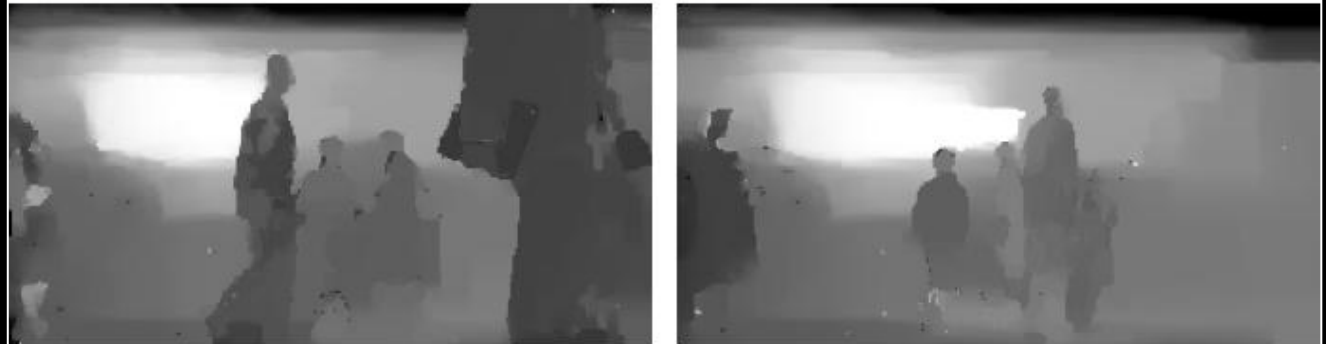


# Application: 2D-to-3D

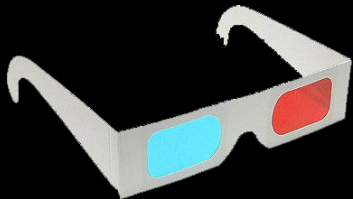
Input



Depth



Anaglyph  
“3D”





# Thanks!

More results, code and dataset available at:  
<http://kevinkarsch.com/depthtransfer>

Our 2D-to-3D

Youtube 2D-to-3D

