# Depth Extraction from Video Using Non-parametric Sampling

Kevin Karsch
University of Illinois

Ce Liu
Microsoft Research
New England

Sing Bing Kang
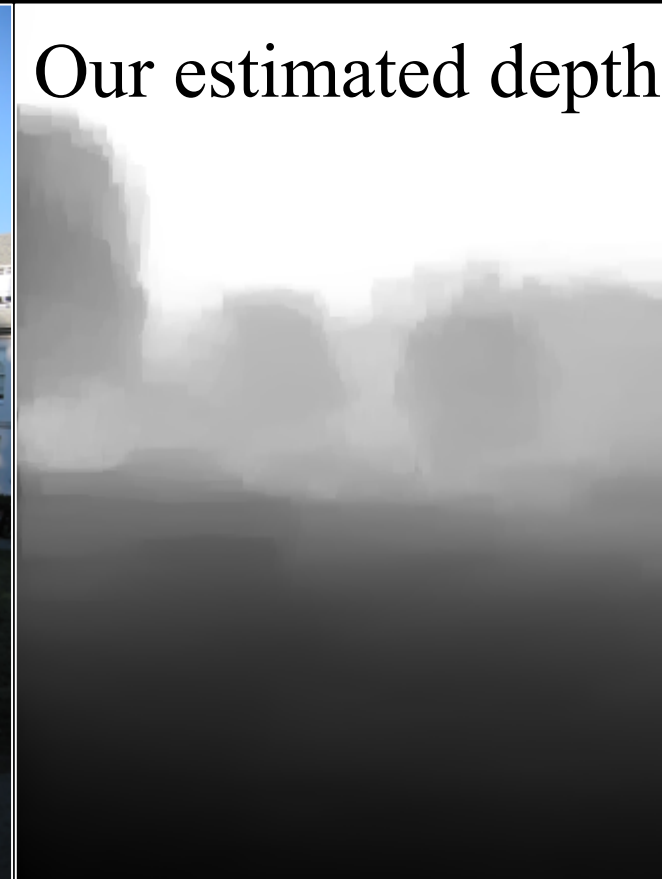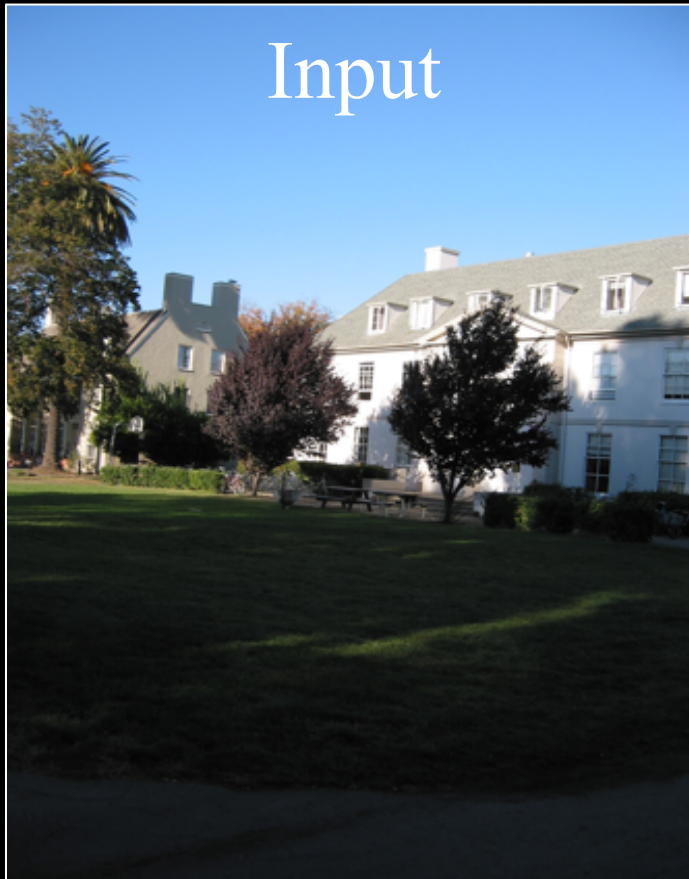Microsoft Research

# 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



Input

Our estimated depth

# Problem Statement

Given an image/video, estimate distance from the camera

No parallax necessary        Camera motion OK        Scene motion OK
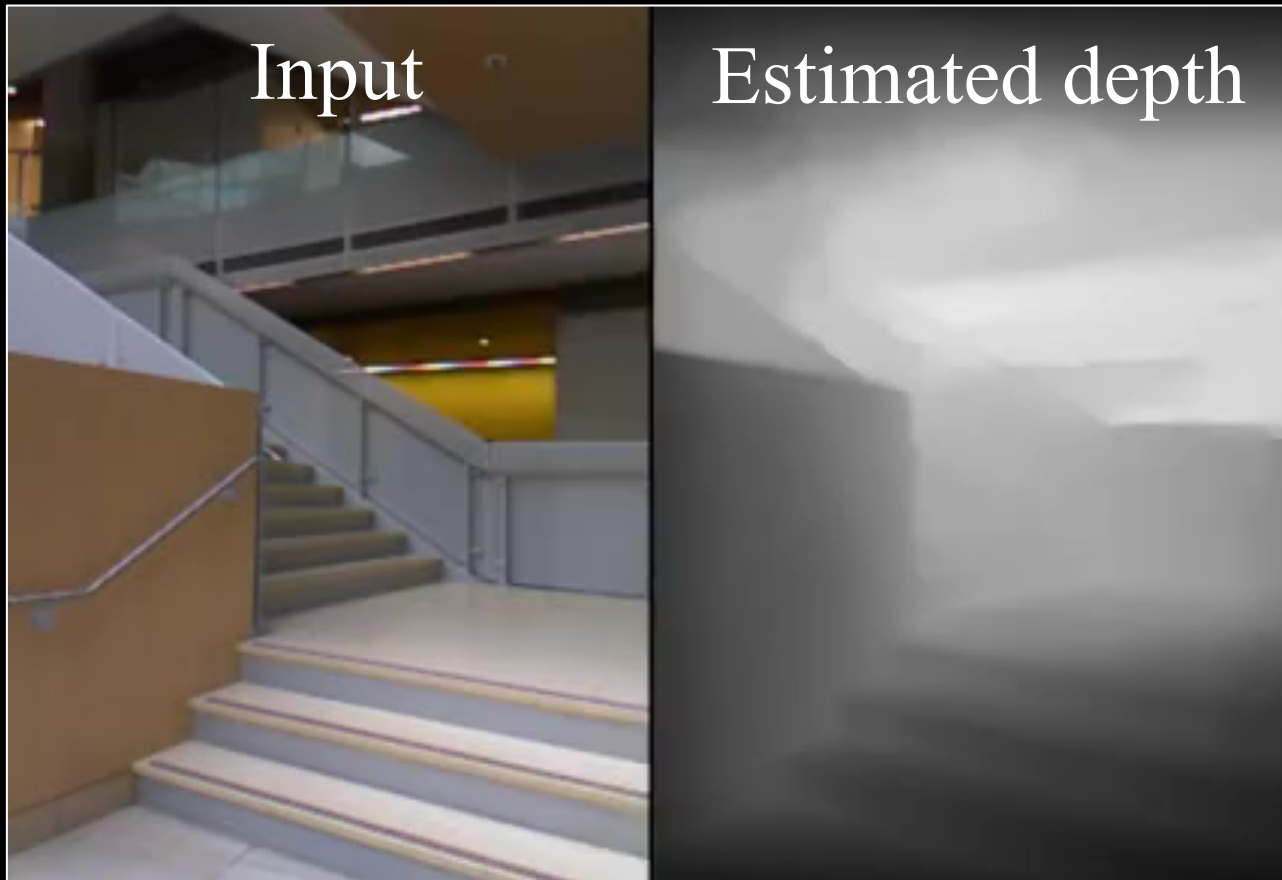
Input        Estimated depth

# Problem Statement

Given an image/video, estimate distance from the camera

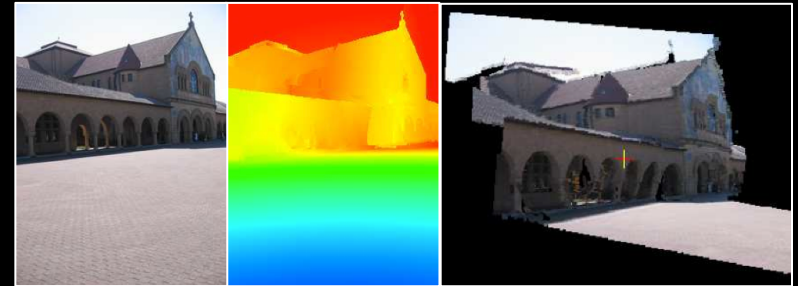No parallax necessary     Camera motion OK     Scene motion OK

# Related Work



[Zhang et al. '09]



[Liu et al. '10]

## Multiview reconstruction

o   Very accurate for videos with moving camera

o   May fail for dynamic scenes

[Newcombe and Davidson '10]

[Furukawa and Ponce '09]

[Zhang et al. '09]

…

## Parametric learning

o   Works well for single images

o   No literature on extending to video

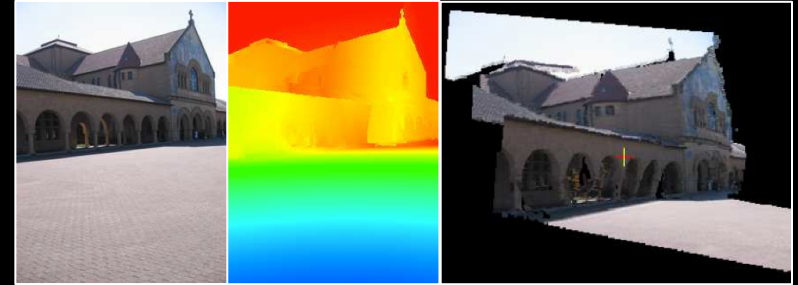[Liu et al. '10]

[Saxena et al. '09]

[Hoiem et al. '05]

…

# Related Work



[Zhang et al. '09]



[Liu et al. '10]

**Multiview reconstruction**

**Parametric learning**

o Very accurate for videos with moving camera

o Works well for single images

o May fail for dynamic scenes

o No literature on extending to video

[Newcombe and Davidson '10]
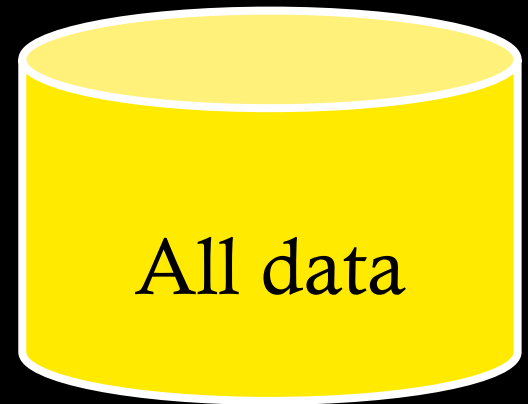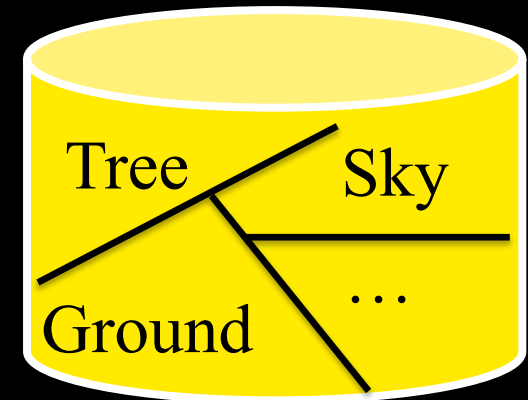[Furukawa and Ponce '09]
[Zhang et al. '09]

…

[Liu et al. '10]
[Saxena et al. '09]
[Hoiem et al. '05]

…

Training set

Pixel level [Saxena et al. '05]

All data

Object level [Lui et al. '10]

Tree　　Sky
Ground　　...

Scene level (ours)

Forest　Office
Street-
level　　...

# RGBD Datasets



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

MSR-V3D
*Indoor scenes*
(Ours)

# SIFT Flow Refresher

o Optical flow using dense SIFT features
  o Larger search window
  o Modified smoothness constraints

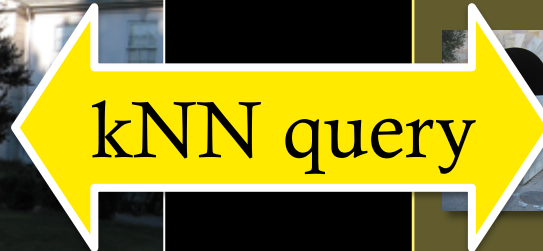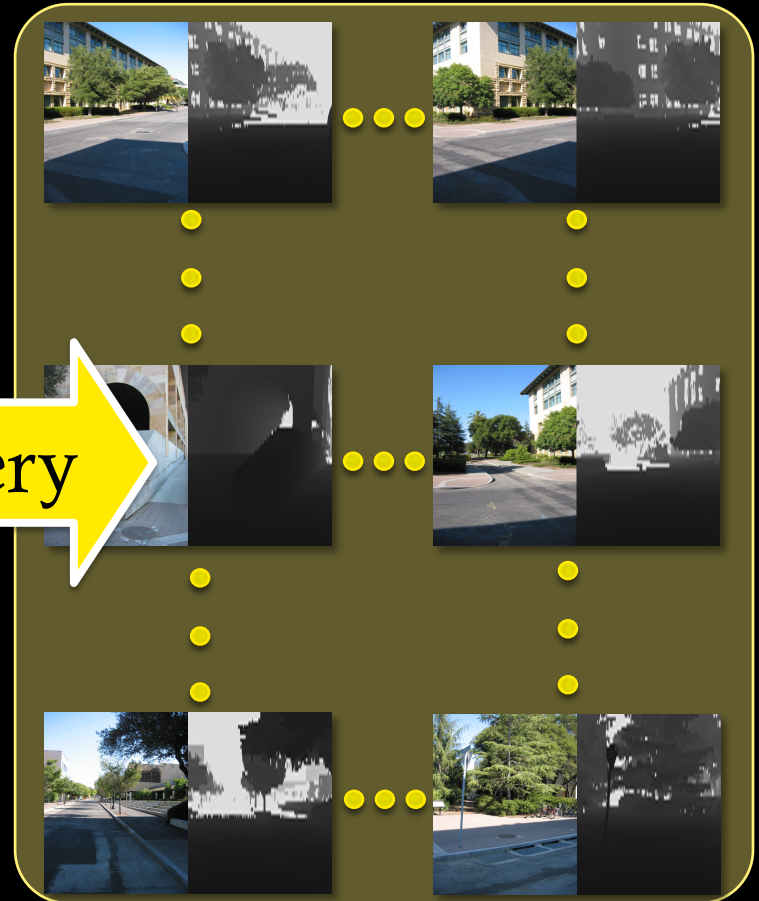o Scenes rearranged so semantics are matched

$\Psi$

Warping
operator

A

B

[Liu et al. '08, '09]

# Algorithm



Input image

RGBD Database

kNN query

# Algorithm

Input

Candidates

Warped candidates



SIFT flow

# Algorithm



Input

Prior

Warped candidates

Depth inference

Inferred depth

# Inference

$$\underset{D}{\arg\min}\, E(D) =$$

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

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

Spatial smoothness

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

$$\underset{D}{\arg\min}\, E(D) =$$

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

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

**Not a discrete MRF!**

Spatial
smoothness

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

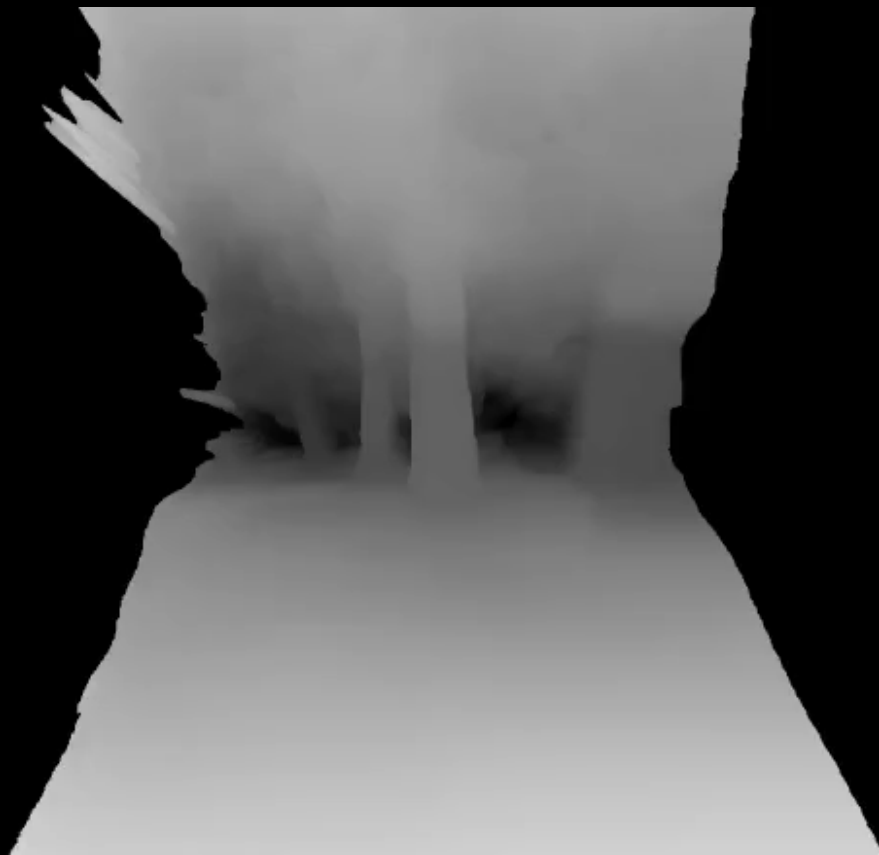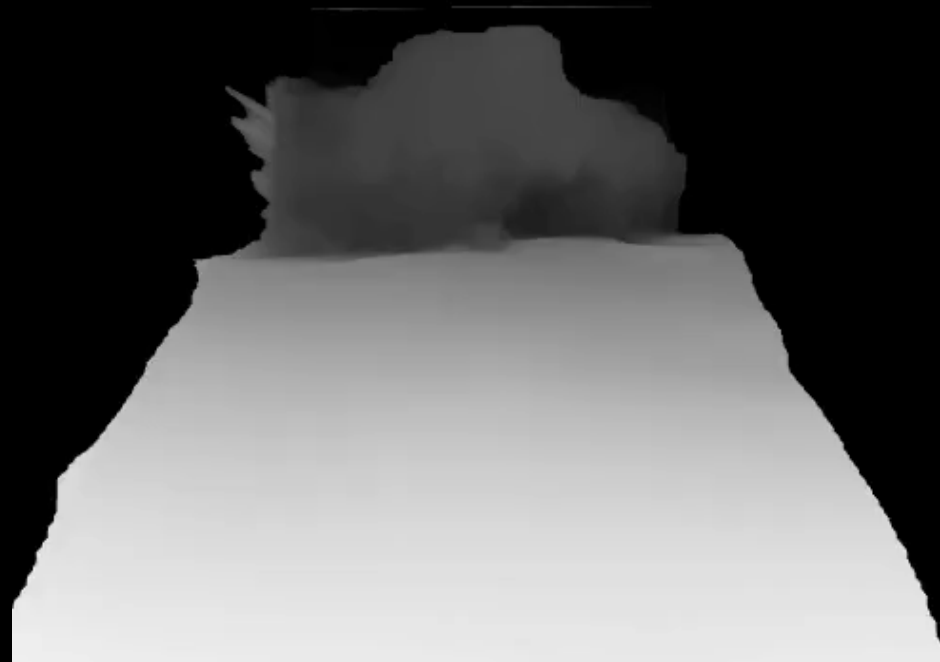Input          True depth          Inferred depth
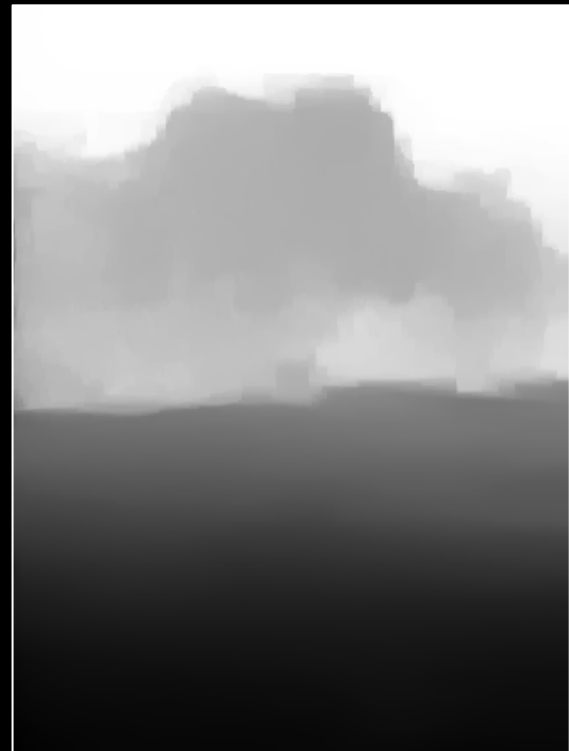
Input          True depth          Inferred depth

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

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

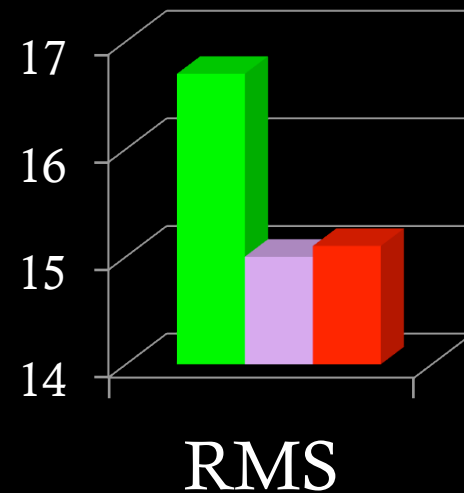Result *without* relative
depth term ($\gamma = 0$)

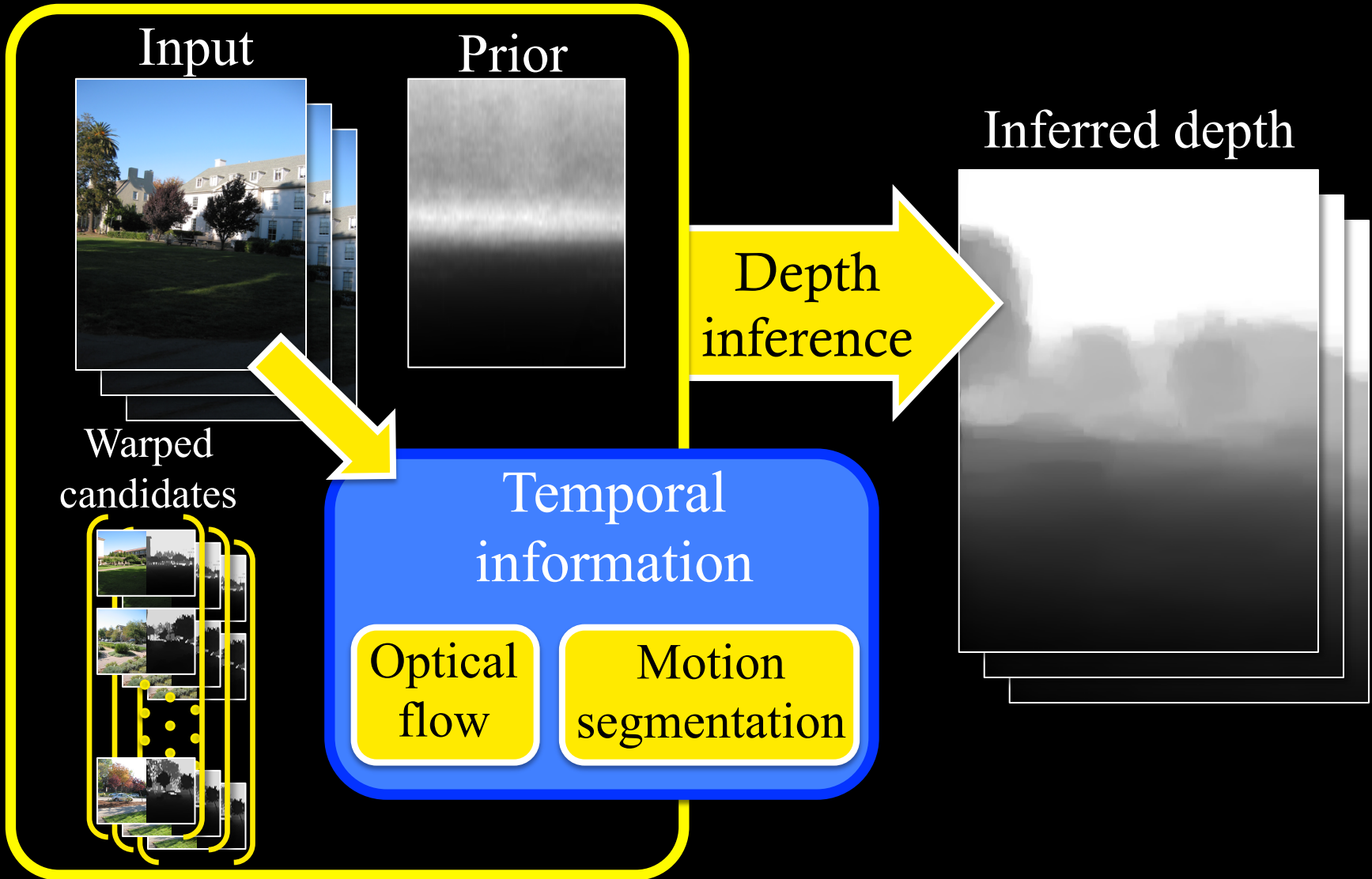Result *with* relative
depth term ($\gamma > 0$)

# 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) |

RGBD+labeled data

RGB | Labels | Depth

# Video Inference

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

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

Single image objective

o    Depth changes are gradual frame-to-frame

o    Moving objects are usually on the ground

# Video Inference

$$\underset{D}{\mathrm{argmin}}\, E_{\mathrm{video}}(D) =$$

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

Single
image
objective

Smooth along
direction of
optical flow

o   Depth changes are gradual frame-to-frame

o   Moving objects are usually on the ground

# Video Inference

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

$$\underset{D}{\arg\min}\, E_{\text{video}}(D) =$$

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

o   Depth changes are gradual frame-to-frame

o   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

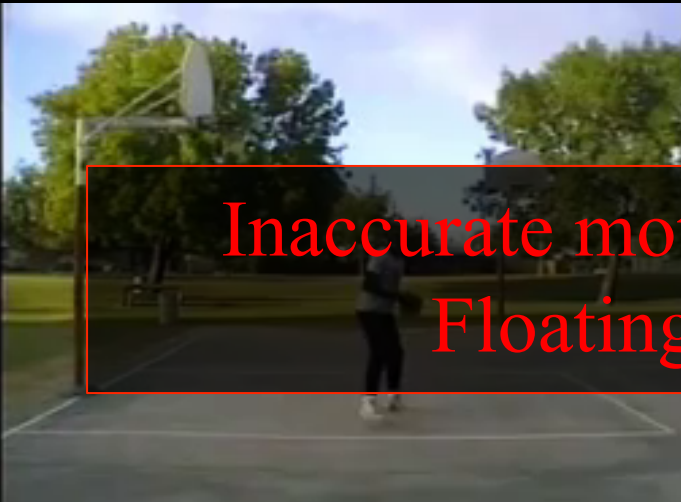Input   Kinect*   Ours   Input   Kinect*   Ours

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

# Limitations



No similar training images
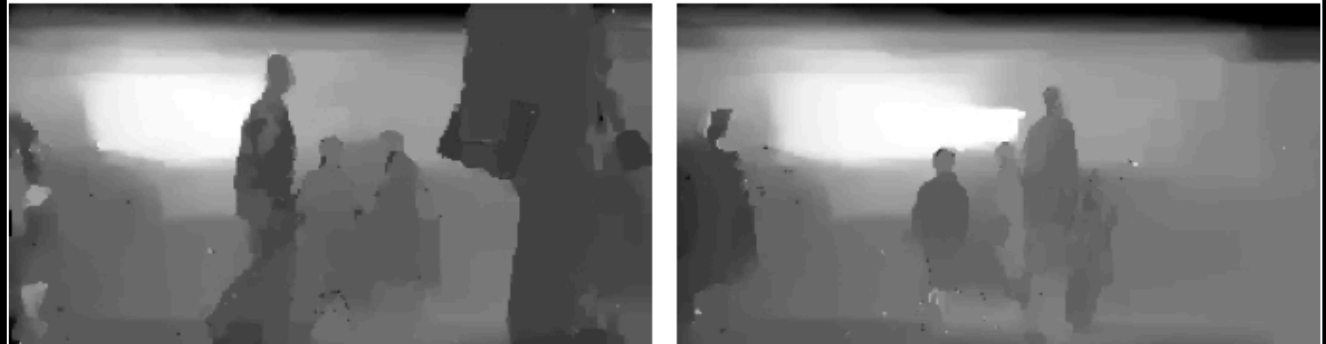
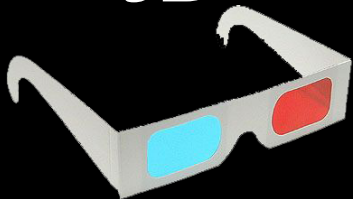Inaccurate motion estimation
Floating objects

# 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