

# 3-D Depth Reconstruction from a Single Still Image

International Journal of Computer Vision (IJCV), Aug 2007

Ashutosh Saxena, Sung H. Chung, Andrew Y. Ng.

April 17, 2013

Presenter: Yiying Li

# Outline

- Motivation
- Previous Attempts
- Methods and Evaluation
  - Feature Extraction
- Evaluation
- Replication
- Difficulties
- Importance

# Motivation

- Recovering depth is an important application in scene understanding, robotics, and 3-d reconstruction.
- Previous work have been extracting depth using stereopsis, or structure from motion.
- Humans use certain monocular cues to indicate depth:
  - Texture variation
  - Gradients
  - Defocus
  - Color and haze

# Previous Attempts

- Reconstruction for known fixed objects (faces, hands), (Nagai et al. 2002)
- Using uniform color and Lambertian surfaces (Many many papers).
- Fourier spectrum to compute mean depth (Torralba and Oliva 2002).
- Supervised learning for 1-D distance for specific obstacles (Michels et al. 2005).
- Fixed sky, ground, vertical regions on the image (Hoeim et al. 2005). No real depth map.

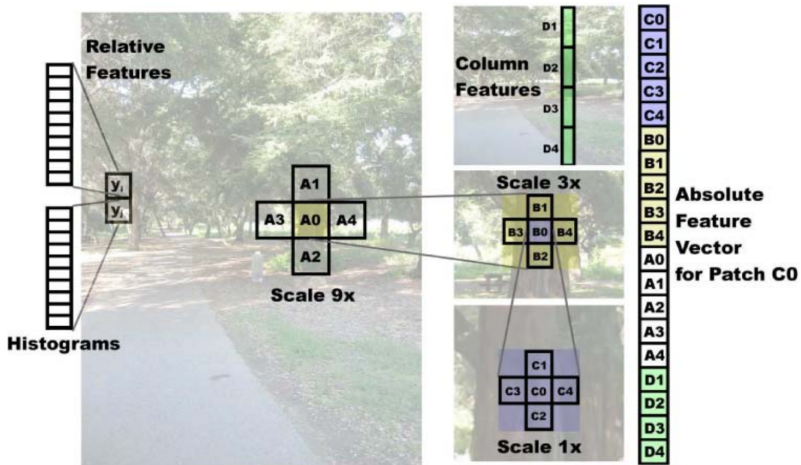
# Feature Extraction

- Monocular cues include: texture variations/gradients, light, haze, defocus, occlusion, known object sizes, and many others.
- If we only look at these features at a local scale we miss out on global properties of the image.
- Convert to YCbCr color space for separation of intensity and color.

# Feature Extraction

- Absolute Features:
  - Texture information is in the intensity channel.
  - Law's mask and six rotated edge filters are convolved with intensity.
  - Haze is usually in low frequency of color channels.
  - Color channels are convolved with a local averaging filter.
  - 17 features: 9 Law's, 6 edge, 2 color.
  - The author also chose to square each filter output as well leaving us with 34 dimensions.
- Relative Features:
  - 10-bin histogram of each of the 17 filter outputs.
  - 170 dimensions.

# Features



# Learning Model

- Represented using a multi-scale hierarchal Markov Random Field (MRF).
- Two different models that only different in inference.
- Gaussian Model:

$$P_G(d|X; \theta, \sigma) = \frac{1}{Z_G} \exp \left( - \sum_{i=1}^M \frac{(d_i(1) - x_i^T \theta_r)^2}{2\sigma_{1r}^2} - \sum_{s=1}^3 \sum_{i=1}^M \sum_{j \in N_s(i)} \frac{(d_i(s) - d_j(s))^2}{2\sigma_{2rs}^2} \right)$$

- Laplacian Model:

$$P_L(d|X; \theta, \lambda) = \frac{1}{Z_L} \exp \left( - \sum_{i=1}^M \frac{(d_i(1) - x_i^T \theta_r)^2}{\lambda_{1r}} - \sum_{s=1}^3 \sum_{i=1}^M \sum_{j \in N_s(i)} \frac{|d_i(s) - d_j(s)|}{\lambda_{2rs}} \right)$$



# Evaluation

- Ground truth depth maps are captured by a 3-D laser scanner.
- Depth maps's are synced to the same FoV of the camera.
- Various different environments, trees, buildings, etc. . .
- 400 training pairs.
- 134 testing pairs.

# Replication Results

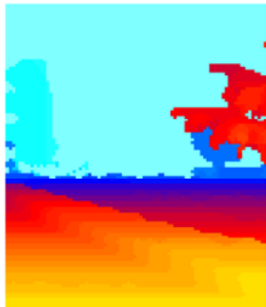
- Replicated using a 8x8 pixel patch in the first scale.
- Author's mean log 10 error for Gaussian Model: 0.133
- Replicated mean log 10 error for Gaussian Model: 0.150
- Replicate mean error: 1.4155 meters

# Compare Replication Results

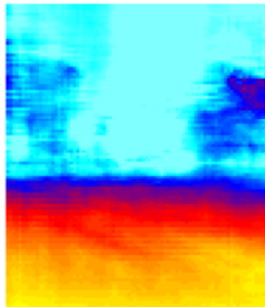
(a) Image



(b) Ground-Truth



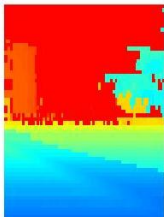
(c) Gaussian



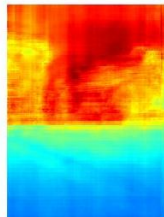
Image



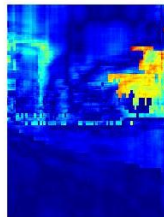
Truth



Generated



Abs Diff

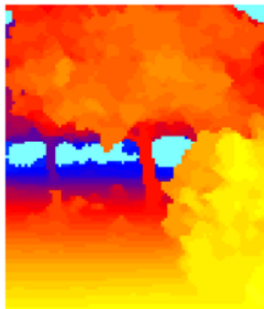


# Compare Replication Results

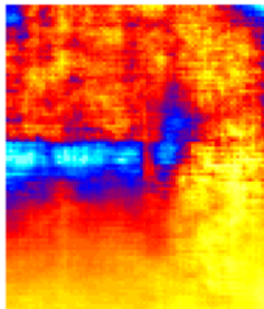
(a) Image



(b) Ground-Truth



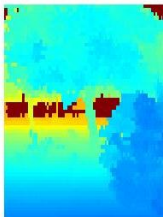
(c) Gaussian



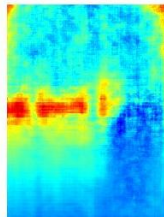
Image



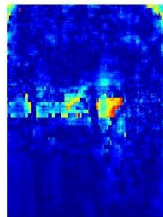
Truth



Generated



Abs Diff

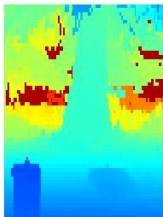


# Good Results

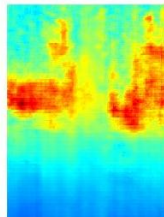
Image



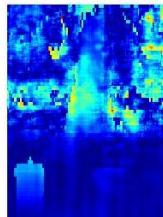
Truth



Generated



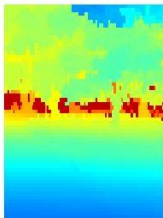
Abs Diff



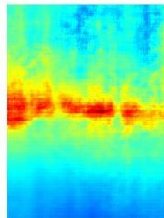
Image



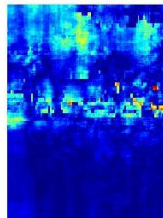
Truth



Generated



Abs Diff

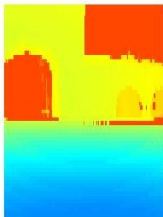


# Good Results

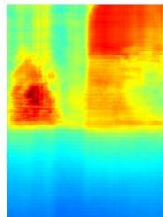
Image



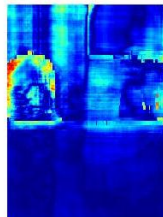
Truth



Generated



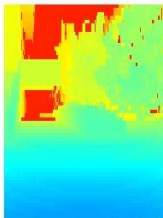
Abs Diff



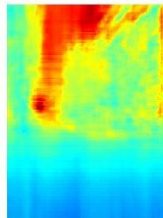
Image



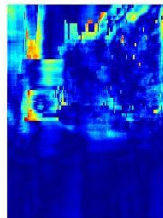
Truth



Generated



Abs Diff

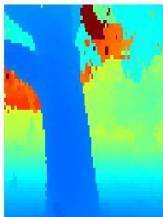


# Bad Results

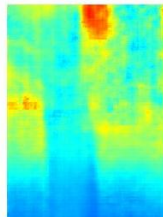
Image



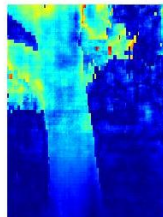
Truth



Generated



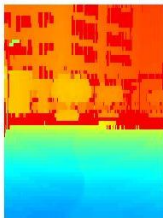
Abs Diff



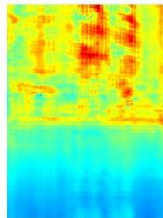
Image



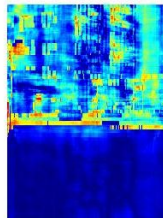
Truth



Generated



Abs Diff

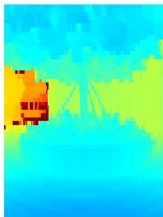


# Bad Results

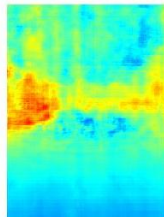
Image



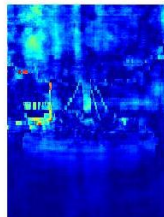
Truth



Generated



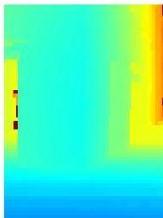
Abs Diff



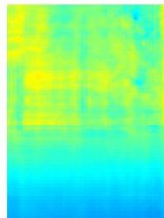
Image



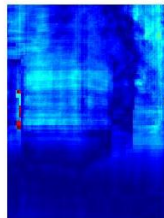
Truth



Generated



Abs Diff





# Difficulties

- Insurmountable difficulties:
  - 2nd term dropped from the model as the paper doesn't provide the MAP inference method.
  - Laplacian wasn't successfully implemented, as I couldn't find a correct way to solve the linear program in the MAP inference.
- Surmountable difficulties:
  - The cost to calculate the absolute features of images
  - 1704x2272 image had a feature "matrix" that cost about 280MB and without vectorization 300s to compute.
  - With some clever optimizations of pre-computing the require information, this was dropped to 30.
  - Can't load all image features to memory (100+ GB) for training.
  - Loading whole image features costs around 300s and around 400s for linear regression.
  - Save each image's features by rows reduces training time from (3 days to 1.3 days).

# Importance and Future Gold

- The problem of generating depth is still an extremely interesting problem that humans seem to be much better at.
- The author takes an interesting approach to generate the right features that have the most information that has a regard to depth.
- The value of the author is more as an augmenting method in order to generate accurate depth maps.
- The work produces a very general depth map that is extremely useful for simple inferences about the environment.
- It is interesting to see what other features detectors we can add to possibly make this method more robust.
- It is a good example just how much information a standard image can hold.
- The method is potentially fast enough given GPU acceleration for real time robotics applications, to augment stereo depth estimation.