

Correspondence Estimation in Images: New Techniques and Applications

Sudipta N. Sinha

Microsoft Research

Correspondence Estimation in Computer Vision

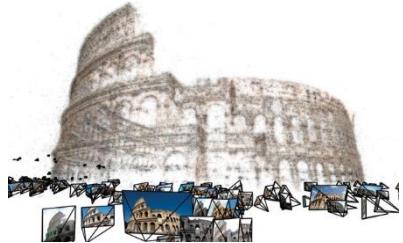


2-view
rigid
sparse

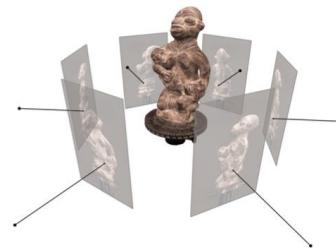
vs.

n-view
non-rigid
dense

Structure from motion



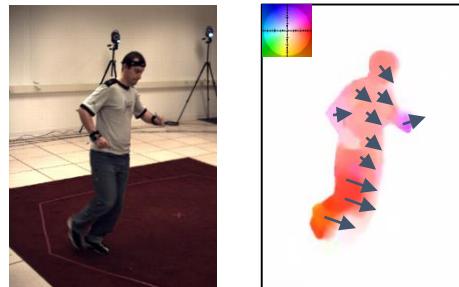
Multiview stereo



Binocular Stereo

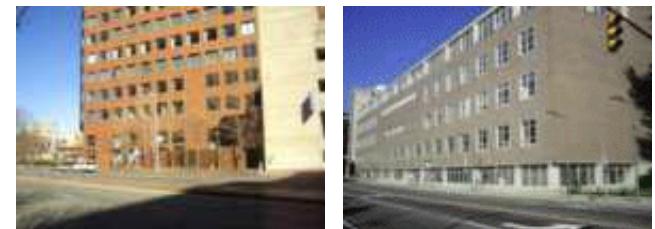


Optical flow



Different scenes

SIFT Flow (Liu et al. 2008)



Deformable Spatial Pyramid Matching (Kim et al. 2013)

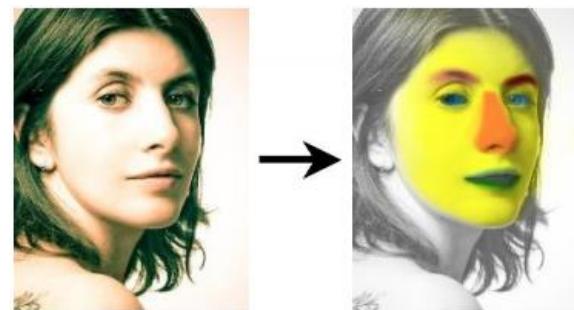


More applications

Label transfer (Face parsing) [Smith et al 2013]



Labeled images



Face skin
Left eye
Right eye
Left brow
Right brow
Nose
Inner mouth
Upper lip
Lower lip
Background

Label Types

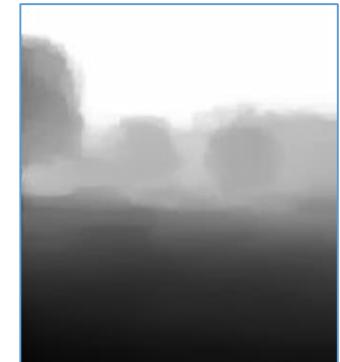
Depth transfer [Karsch et al. 2012]



RGB-D database



Query



Predicted Depth

Overview

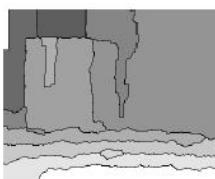
- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Surface-based stereo

- Piecewise planar stereo



Birchfield and Tomasi 2001

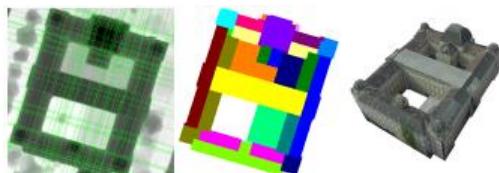


Furukawa et al. 2008



Sinha et al. 2009

- Surface stereo



Zebедин et al. 2008



Gallup et al. 2010

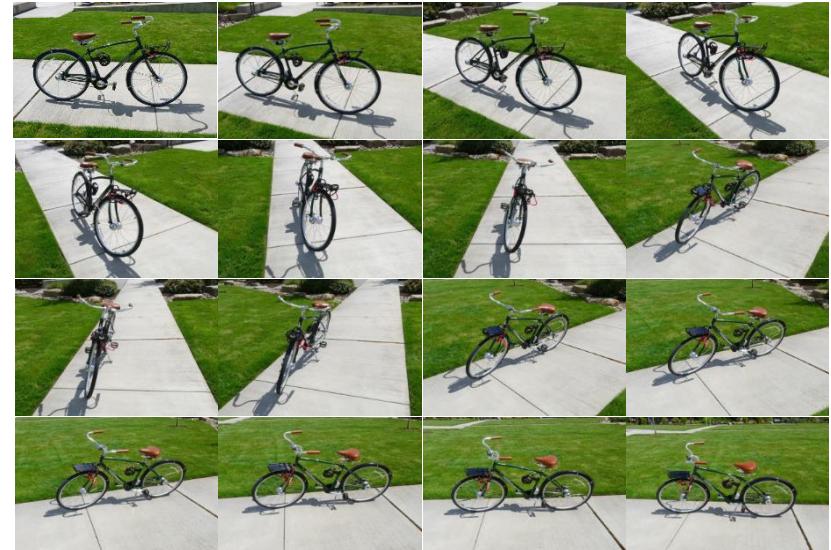


Bleyer et al. 2010, 2011

Multiple View Object Cosegmentation using Appearance and Stereo Cues

Kowdle, Sinha and Szeliski (ECCV 2012)

Input
Images



Final
Results

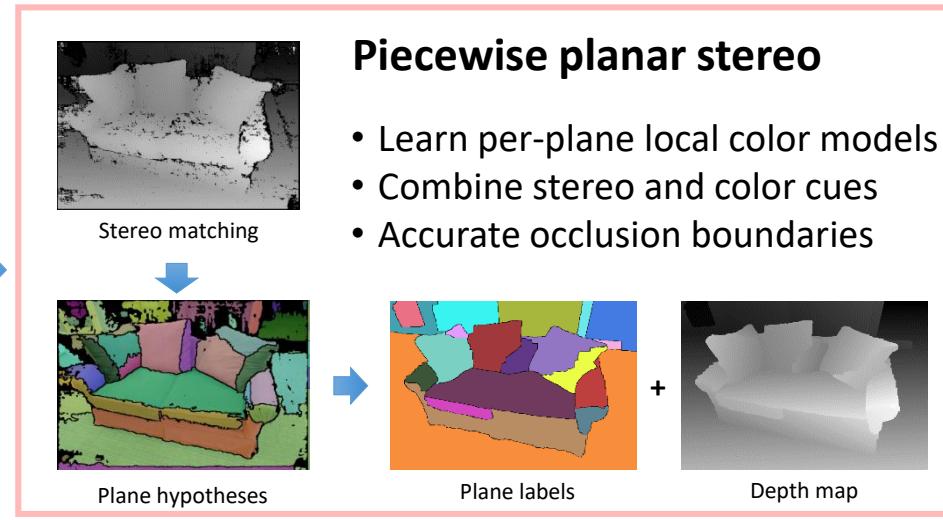


Multiple View Object Cosegmentation using Appearance and Stereo Cues

Kowdle, Sinha and Szeliski (ECCV 2012)



Input images

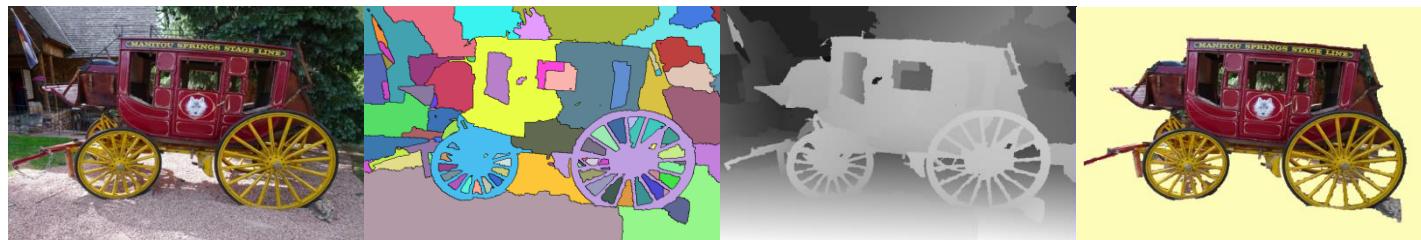
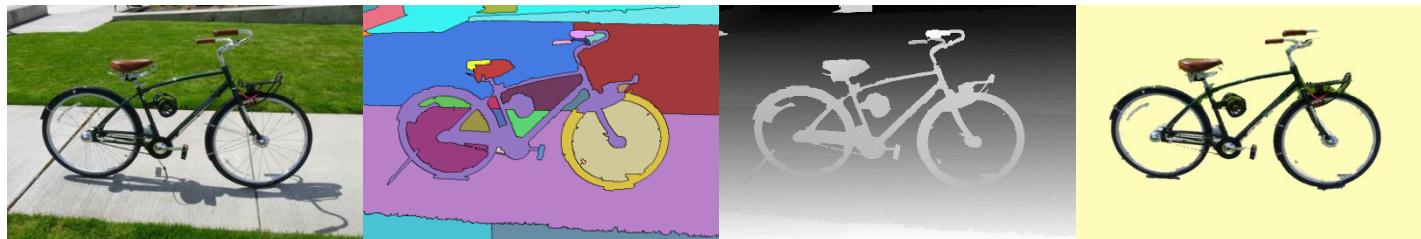


Multiview Foreground Object Segmentation

- Infer what constitutes the foreground *object*
- Soft segmentation consistency in multiple-views

Multiple View Object Cosegmentation using Appearance and Stereo Cues

Kowdle, Sinha and Szeliski (ECCV 2012)



Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - **High Resolution Stereo Matching**
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Stereo benchmarks

Middlebury
(v2 now offline)

[vision.middlebury.edu](http://vision.middlebury.edu/stereo/)

Stereo Evaluation Datasets Code Submit

Middlebury Stereo Evaluation - Version 2

New features and main differences to version 1. [Submit and evaluate your own results](#)

Open a new window for each link

| Error Threshold = 1 | | Sort by nonocc | | | Sort by all | | | Sort by disc | | | Average percent of bad pixels (explanation) |
|---------------------|------|--|--|--|--|-----|------|--------------|-----|------|---|
| Error Threshold... | | nonocc | all | disc | nonocc | all | disc | nonocc | all | disc | |
| Algorithm | Avg. | Tsukuba ground truth | Venus ground truth | Teddy ground truth | Cones ground truth | | | | | | |
| ADCensus [82] | 10.9 | 1.07 ₁₆ 1.48 ₁₃ 5.73 ₁₉ | 0.09 ₂ 0.25 ₉ 1.15 ₂ | 4.10 ₁₃ 6.22 ₅ 10.9 ₁₁ | 2.42 ₁₄ 7.25 ₁₁ 6.95 ₁₅ | | | | | | 3.97 |
| AdaptingBP [16] | 14.2 | 1.11 ₁₉ 1.37 ₈ 5.79 ₂₁ | 0.10 ₄ 0.21 ₈ 1.44 ₁₀ | 4.22 ₁₅ 7.06 ₁₂ 11.8 ₁₆ | 2.48 ₁₈ 7.92 ₂₃ 7.32 ₂₃ | | | | | | 4.23 |
| CoopRegion [39] | 14.8 | 0.87 ₄ 1.16 ₁ 4.61 ₄ | 0.11 ₅ 0.21 ₅ 1.54 ₁₀ | 5.18 ₂₅ 8.31 ₁₆ 13.0 ₂₁ | 2.79 ₃₅ 7.18 ₁₀ 8.01 ₄₀ | | | | | | 4.41 |
| RDP [87] | 19.2 | 0.97 ₉ 1.39 ₁₀ 5.00 ₈ | 0.21 ₃₄ 0.38 ₂₄ 1.89 ₁₉ | 4.84 ₁₈ 9.94 ₂₅ 12.6 ₁₈ | 2.53 ₂₁ 7.69 ₁₆ 7.38 ₂₄ | | | | | | 4.57 |
| MultiRBF [129] | 19.6 | 1.33 ₄₁ 1.56 ₁₇ 6.02 ₂₈ | 0.13 ₈ 0.17 ₂ 1.84 ₁₆ | 5.09 ₂₄ 6.36 ₇ 13.4 ₂₈ | 2.90 ₄₂ 6.76 ₅ 7.10 ₂₀ | | | | | | 4.39 |
| DoubleBP [34] | 20.0 | 0.88 ₆ 1.28 ₉ 4.76 ₇ | 0.13 ₉ 0.45 ₄₁ 1.87 ₁₈ | 3.53 ₁₀ 8.30 ₁₆ 9.83 ₁₃ | 2.90 ₄₁ 8.78 ₅₀ 7.79 ₃₂ | | | | | | 4.19 |
| MDPM [140] | 20.3 | 1.15 ₂₀ 1.59 ₂₀ 6.14 ₃₁ | 0.14 ₁₄ 0.36 ₂₂ 1.52 ₉ | 3.79 ₁₁ 5.78 ₄ 11.1 ₁₃ | 2.74 ₂₉ 8.38 ₃₈ 9.71 ₃₅ | | | | | | 4.22 |
| OutlierConf [40] | 20.5 | 0.88 ₅ 1.43 ₁₂ 4.74 ₆ | 0.18 ₂₃ 0.26 ₁₁ 2.40 ₃₄ | 5.01 ₂₀ 9.12 ₂₄ 12.8 ₂₀ | 2.78 ₃₄ 8.57 ₄₁ 6.99 ₁₈ | | | | | | 4.60 |
| AdaptiveGF [127] | 24.1 | 1.04 ₁₂ 1.53 ₁₄ 5.62 ₁₄ | 0.17 ₂₂ 0.41 ₃₂ 1.98 ₂₂ | 5.71 ₃₅ 11.3 ₄₁ 14.3 ₃₂ | 2.44 ₁₈ 8.22 ₃₀ 7.05 ₁₈ | | | | | | 4.98 |

450 x 376 pixels
D ≈ 16...60

KITTI

The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago



home setup stereo flow odometry object tracking road raw data submit results jobs
Andreas Geiger (MPI Tübingen) | Philip Lenz (KIT) | Christoph Stiller (KIT) | Raquel Urtasun (University of Toronto)

Stereo Evaluation

| Rank | Method | Setting | Code | Out-Noc | Out-All | Avg-Noc | Avg-All | Density | Runtime | Environment | Compare |
|--|------------|---------|------|---------|---------|---------|---------|----------|---------|------------------------------------|---------|
| 1 | SceneFlow | | | 2.98 % | 3.97 % | 0.8 px | 1.0 px | 100.00 % | 35 s | 1 core @ 3.5 Ghz (C/C++) | |
| Anonymous submission | | | | | | | | | | | |
| 2 | PCBP-SS | | | 3.40 % | 4.72 % | 0.8 px | 1.0 px | 100.00 % | 5 min | 4 cores @ 2.5 Ghz (Matlab + C/C++) | |
| K. Yamaguchi, D. McAllister and R. Urtasun: Robust Monocular Epipolar Flow Estimation CVPR 2013. | | | | | | | | | | | |
| 3 | gtRF-SS | | | 3.83 % | 4.59 % | 0.9 px | 1.0 px | 100.00 % | 1 min | 1 core @ 2.5 Ghz (Matlab + C/C++) | |
| Anonymous submission | | | | | | | | | | | |
| 4 | StereoSLIC | | | 3.92 % | 5.11 % | 0.9 px | 1.0 px | 99.89 % | 2.3 s | 1 core @ 3.0 Ghz (C/C++) | |
| K. Yamaguchi, D. McAllister and R. Urtasun: Robust Monocular Epipolar Flow Estimation CVPR 2013. | | | | | | | | | | | |
| 5 | PR-SfE | | | 4.02 % | 4.87 % | 0.9 px | 1.0 px | 100.00 % | 200 s | 4 cores @ 3.0 Ghz (Matlab + C/C++) | |
| C. Vogel, S. Roth and K. Schindler: Piecewise Rigid Scene Flow International Conference on Computer Vision (ICCV) 2013. | | | | | | | | | | | |
| 6 | PCBP | | | 4.04 % | 5.37 % | 0.9 px | 1.1 px | 100.00 % | 5 min | 4 cores @ 2.5 Ghz (Matlab + C/C++) | |
| K. Yamaguchi, T. Urtasun, D. McAllister and R. Urtasun: Continuous Monocular Depth Map for Robust Stereo Estimation ECCV 2012. | | | | | | | | | | | |

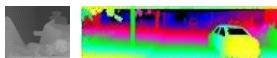
1241 x 376 pixels
D ≈ 70...150

Kim et al. 2013 (Disney Research)

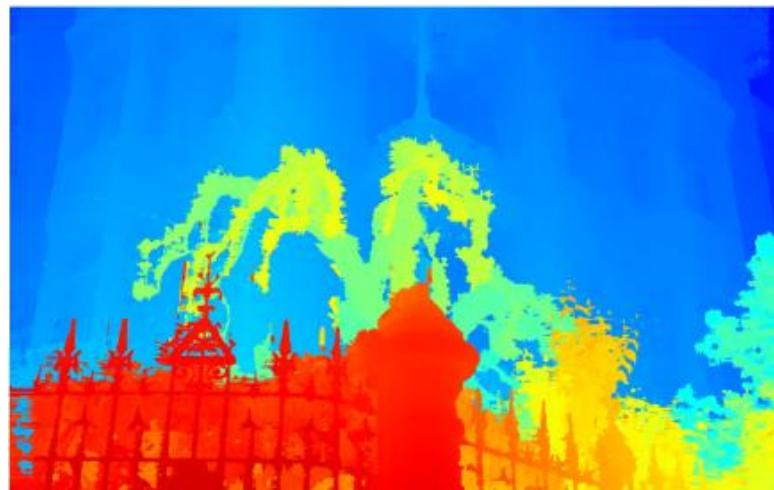
Middlebury
Cones/Teddy

Middlebury v3

KITTI



5 – 6 MPixels



~20 MPixels

Disparity Search Space

| Middlebury (old) | KITTI |
|------------------|----------------|
| 10 Mdisp. | 40 Mdisp. |
| Middlebury New | Disney Mansion |
| 1.5 Gdisp. | 20 Gdisp. |

$O(P*D) = O(s^3)$ for resolution s , P : pixels,

D : disparities

(Most) methods are $O(P*D)$, or $O(P*D^2)$; they do not scale

Our Goals:

- Ideally, want $O(P)$
- Avoid enumerating all disparities
- Optimization should scale

Related Work

- Efficient approximate energy minimization
 - Semi-global Matching (SGM) [Hirschmüller 2005]
- Disparity Refinement [Ma 2013, ...]
- Avoid exploring the whole DSI
 - Coarse-to-fine [long tradition]
 - Seed & Grow [Cech & Sara 2007, ...]
 - PatchMatch stereo [Bleyer et al. 2011]
 - ELAS [Geiger et al. 2010]
 - Bilateral space edge-aware stereo [Barron et al. 2015]
 - Tunable Stereo [Pillai et al. 2016]

Efficient High-Resolution Stereo Matching using Local Plane Sweeps

CVPR 2014



Sudipta Sinha

Microsoft Research

Daniel Scharstein

Middlebury College

Richard Szeliski*

Facebook



* while at Microsoft Research

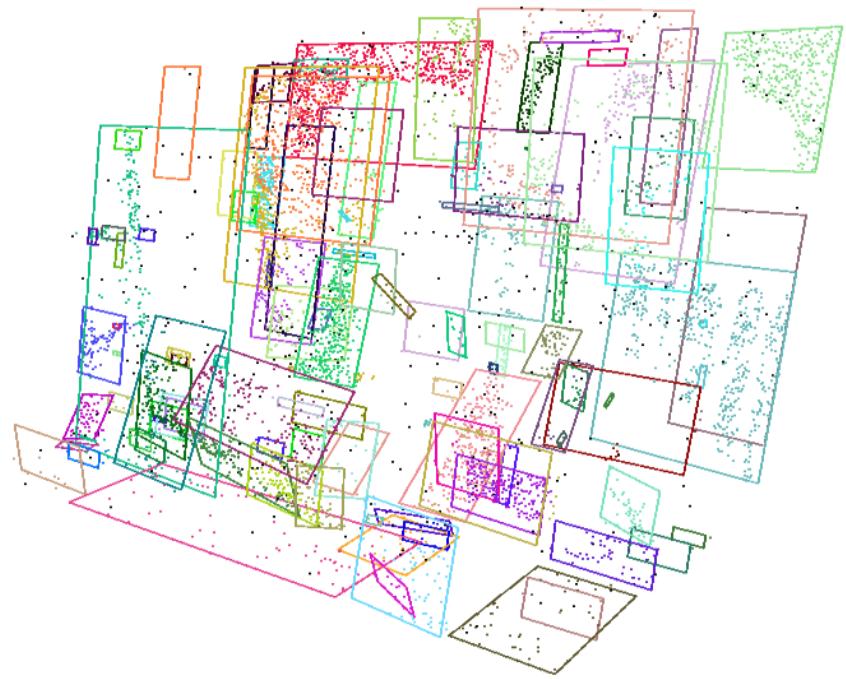
Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)

Plane Hypothesis Generation



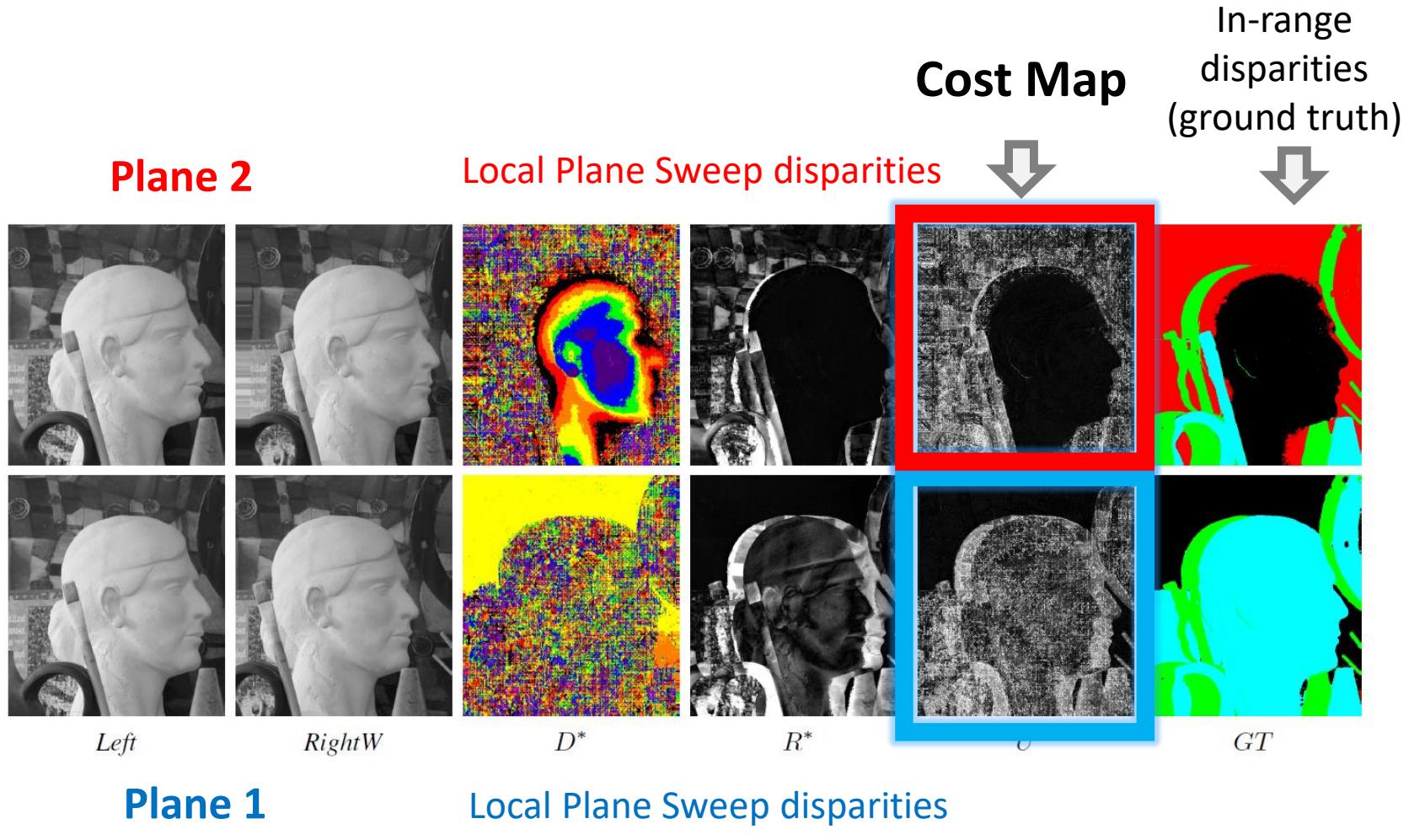
Spurious Classifications



Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around planes
 - narrow disparity range; SGM optimization

Local Plane Sweeps



Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around planes
 - narrow disparity range; SGM optimization

Impose Tile structure

- Perform LPS on tiles and propagate planes to adjoining tiles

Local Plane Sweep Stereo

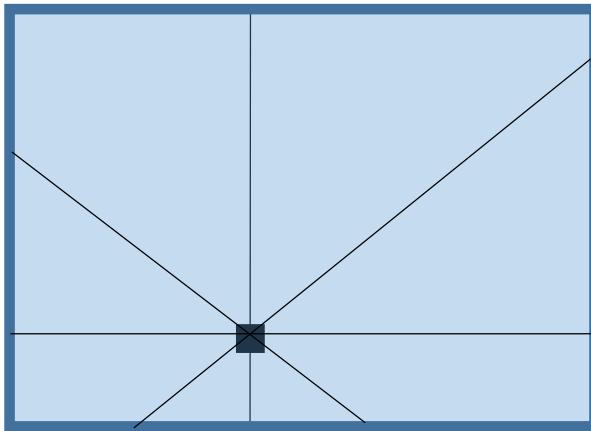
- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around planes
 - narrow disparity range; SGM optimization

Impose Tile structure

- Perform LPS on tiles and propagate planes to adjoining tiles
- Global optimization
 - Assign pixels to surface proposals
 - Approximate energy minimization (via SGM)
 - Extend SGM to exploit tile structure and sparse label sets

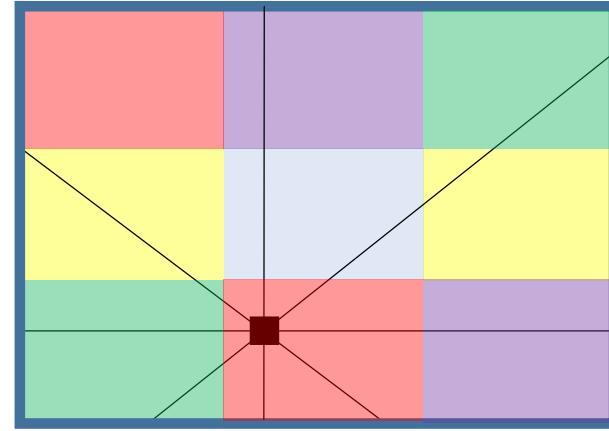
Global Optimization (via SGM)

- Message passing on 1D paths (8 directions)



SGM

- same labels at all pixels



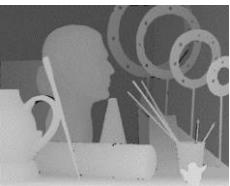
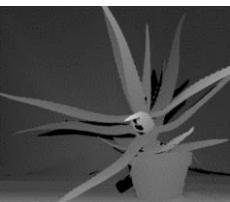
LPS – SGM

- Label sets vary tile to tile*
- Needs book-keeping at tile boundaries*

Datasets

Midd9

2003-2006 Middlebury
(1.4 – 2.7 MP)



New7

2011-2014 Middlebury
(5.1 – 6.0 MP)



Disney

Kim et al. 2013, in SIGGRAPH
(4.5 – 20 MP)

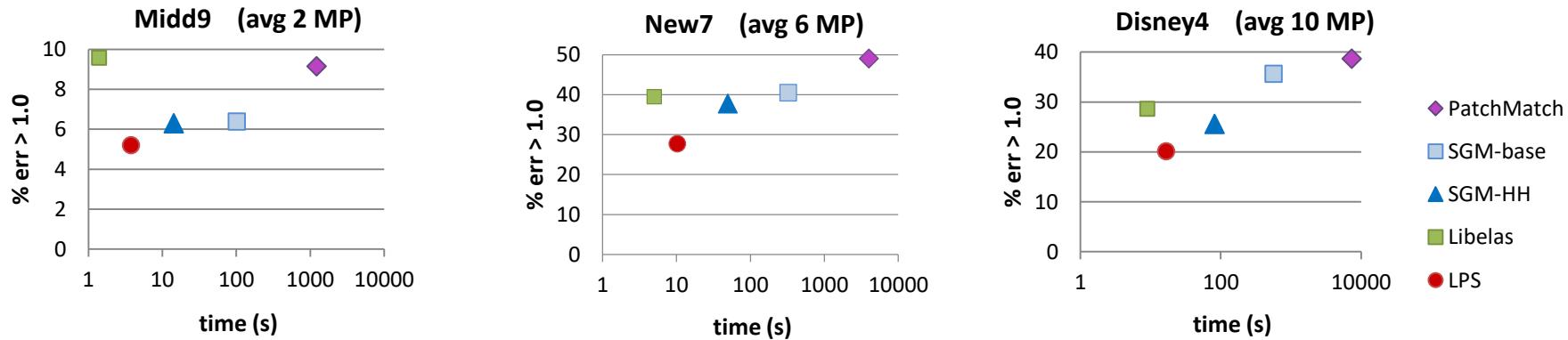


Experiments

- Evaluation:
 - - PatchMatch Stereo [Bleyer et al. 2011]
 - - SGM (our impl.)
 - - SGM-HH [Hirschmüller 2005]
 - - ELAS [Geiger et al. 2010]
 - - LPS
- Metric:
 - 1 pixel disparity error at non-occluded pixels.

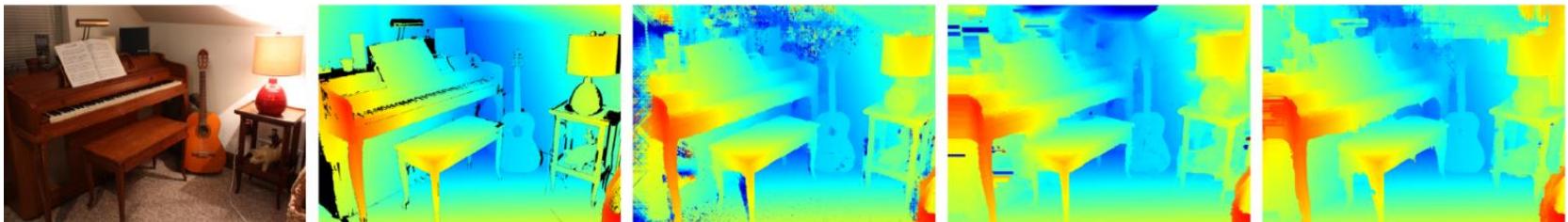
Results – Accuracy vs. Runtime

Error vs. runtime, 1.0 pixel threshold

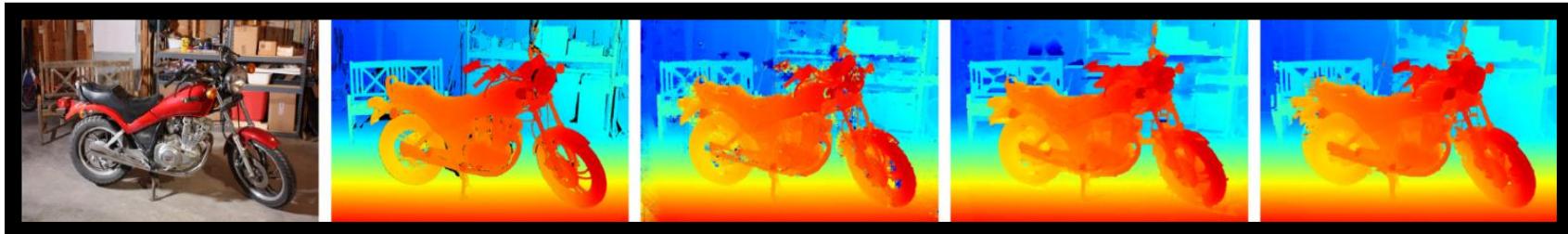


- LPS is the most accurate followed by SGM
- ELAS is the fastest, LPS is 2nd.
 - no GPUs were used.

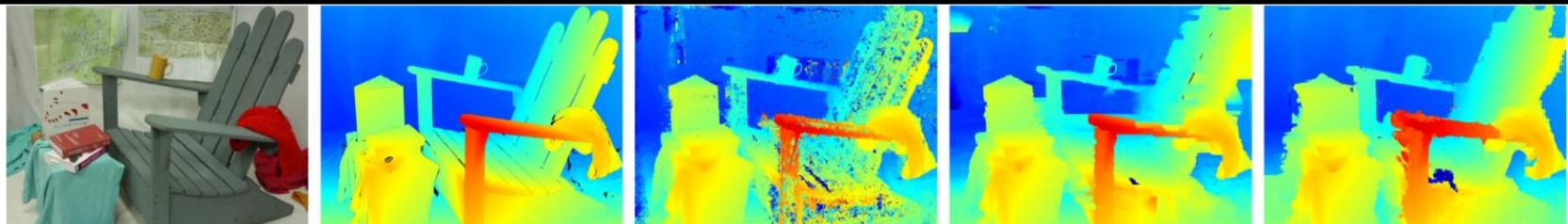
Piano



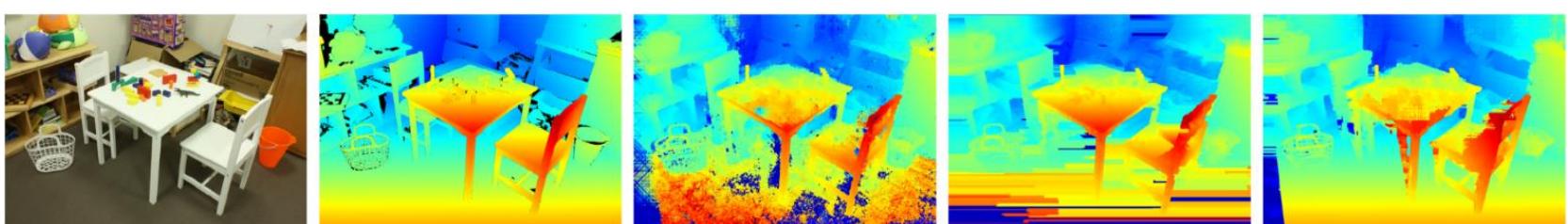
Motorcycle



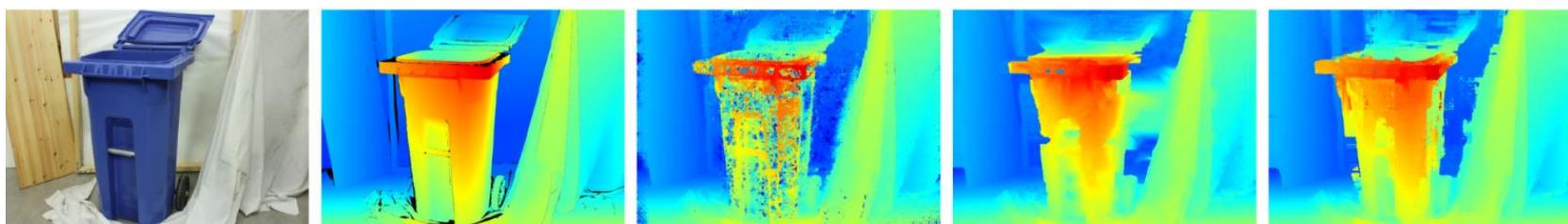
adirondack



Playable



Recycle



Left Image

Ground Truth

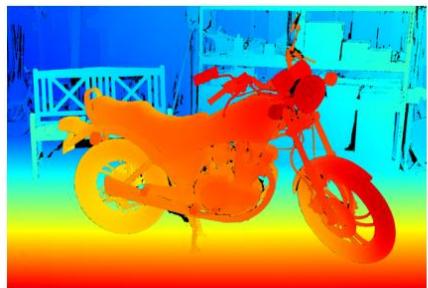
SGM

ELAS

LPS (ours)

Motorcycle (1 pixel error maps)

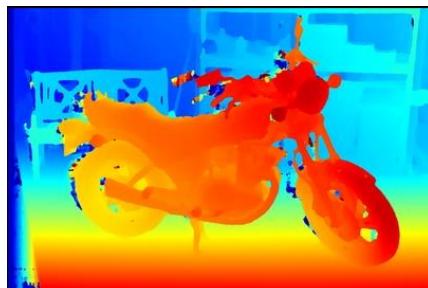
Ground truth



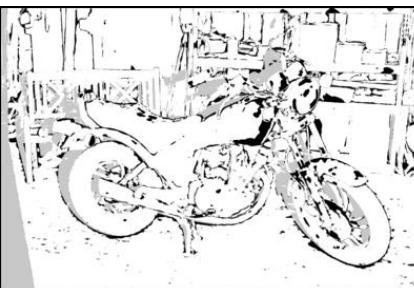
Occlusion mask



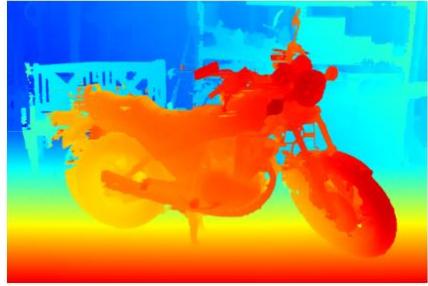
MC-CNN-acrt (109 seconds)



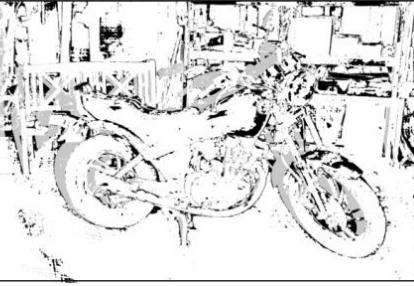
err1 = 9.6



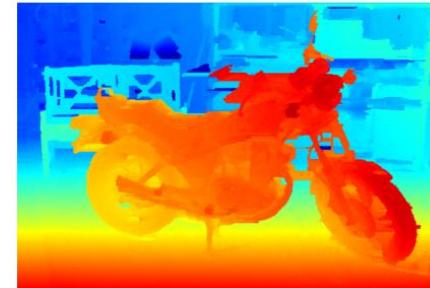
LPS (9.6 seconds)



err1 = 12.2



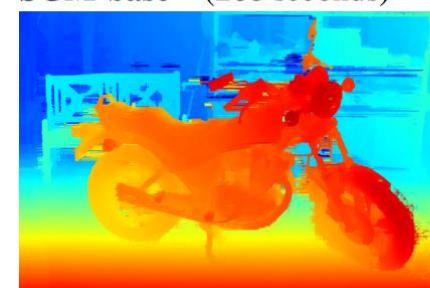
Libelas (5.0 seconds)



err1 = 34.0



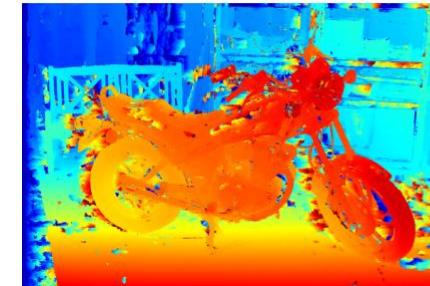
SGM-base (268 seconds)



err1 = 34.4



PatchMatch (3330 seconds)



err1 = 33.8



Summary

Advantages

- Avoid exploring full search space
- Runtime independent of disparity range
- Handles weakly textured slanted surfaces

Limitations

- Can miss surfaces not among initial proposals
- No good “stopping criterion” for proposal generation
- Unclear how to incorporate coarse to fine reasoning

Promising Directions

- Avoid monolithic optimization
- Residual analysis to guide efficient search

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- **Joint Correspondence and Cosegmentation**
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Joint Cosegmentation and Dense Correspondence Estimation

CVPR 2016 (to appear)



Tatsunori Taniai

Univ. of Tokyo



Yoichi Sato

Univ. of Tokyo

Sudipta Sinha

Microsoft Research

Problem

Input

Image pair containing semantically related objects but different instances



Source



Target

Problem

Input

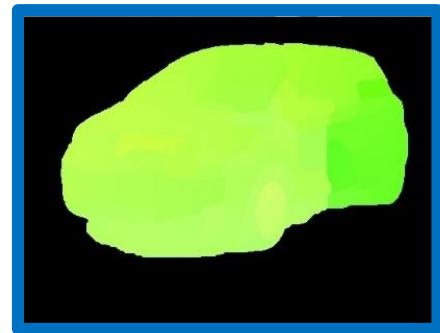
Image pair containing semantically related objects but different instances



Source



Target



Mask + Flow



Warped Source
Image

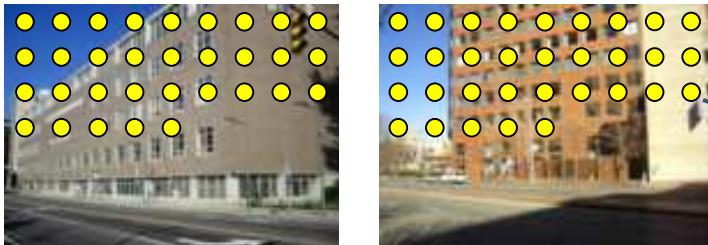
Output

Find the common region i.e. ***foreground (binary) mask***

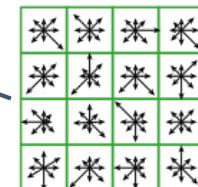
Find the dense ***flow map*** associated with *foreground*.

SIFT Flow [Liu+ 2008]

Robust visual-similarity matching using dense SIFT

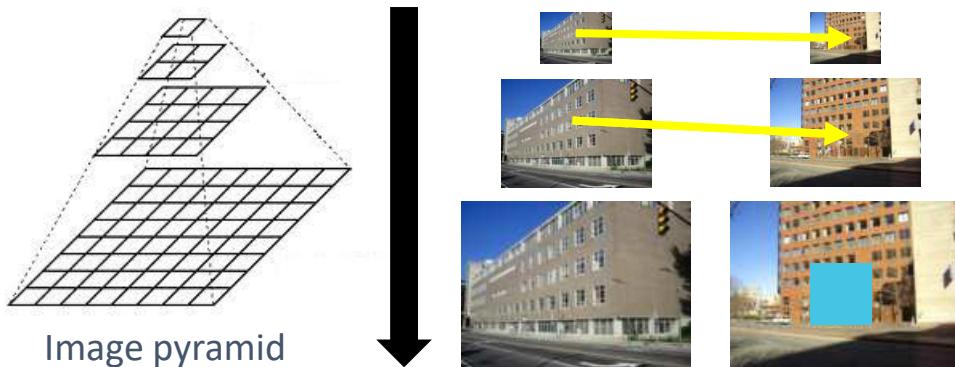


Dense SIFT descriptors



- At every pixel
- Same scale
- No rotation

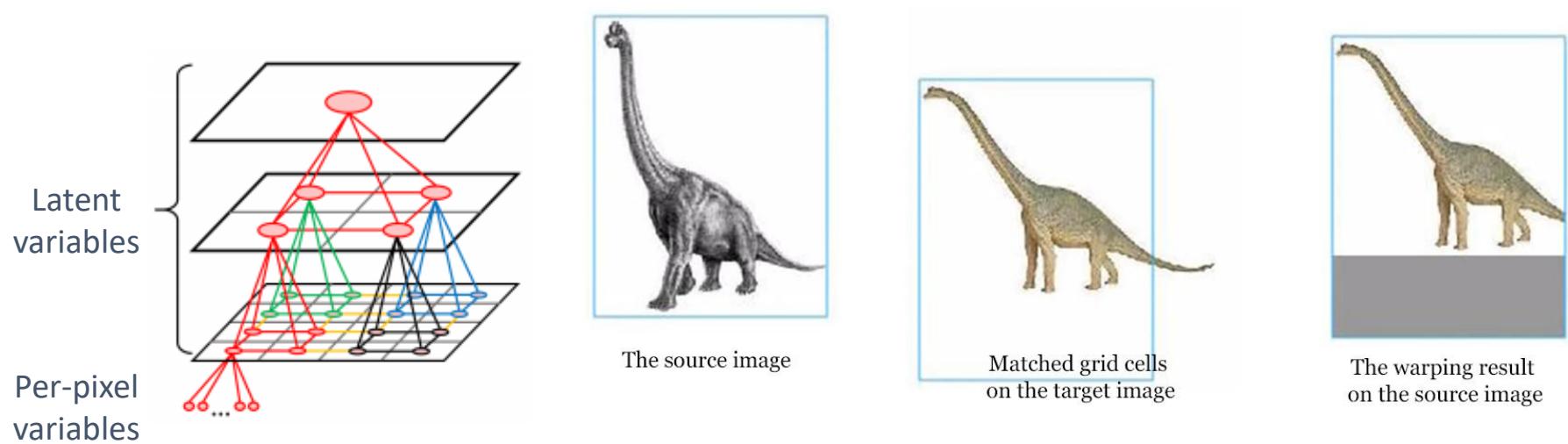
Efficient coarse to fine inference [Felzenszwalb and Huttenlocher 2004]



- Estimate and propagate from coarse to fine levels
- Search in a limited range from propagated points

Generalized Deformable Spatial Pyramids

[Hur+15, Kim+13]



- Powerful yet flexible regularization
- Hierarchy is not segmentation/flow aware

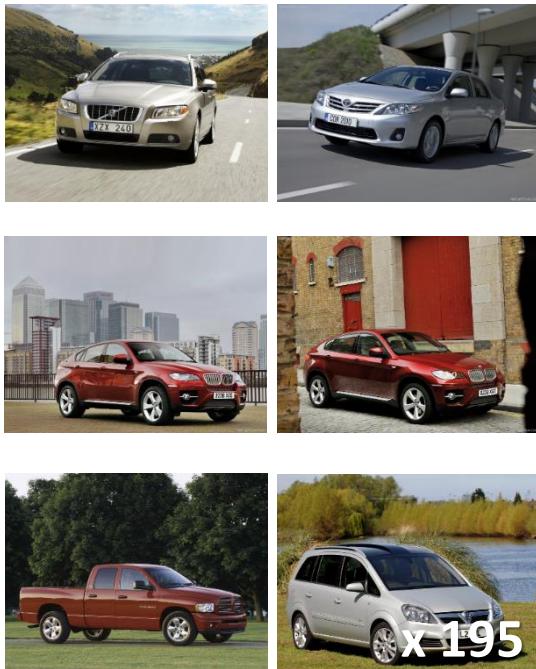
Image Co-segmentation

- Common region shares similar statistics
- Pixel correspondence in common region not modeled



New Dataset

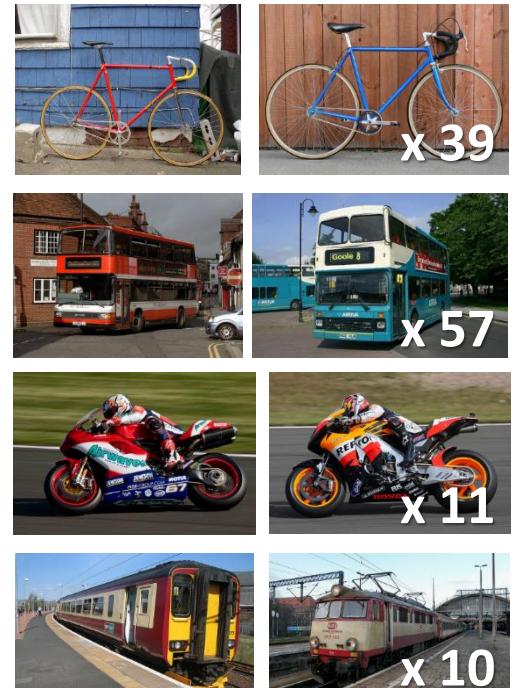
FG3DCar ^[1]



JODS ^[2]



PASCAL ^[3]



[1] Lin et al 2014, "Jointly Optimizing 3D Model Fitting and Fine-Grained Classification"

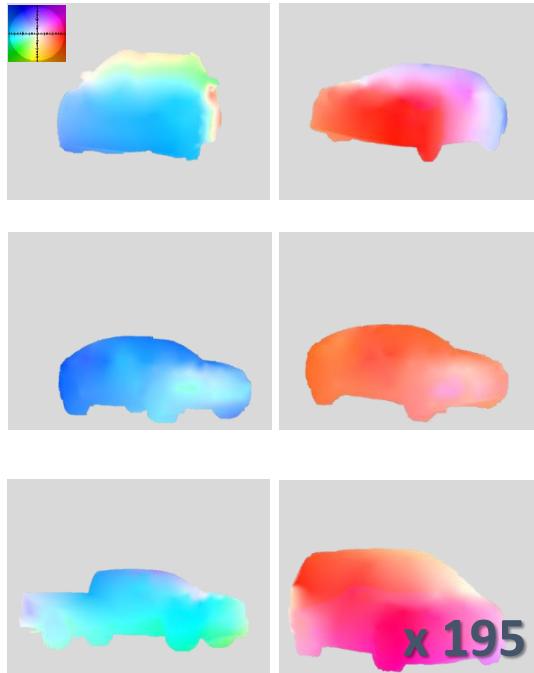
[2] Rubinstein et al 2013 "Unsupervised Joint Object Discovery and Segmentation in Internet Images"

[3] Hariharan et al 2011, "PASCAL segmentation dataset"

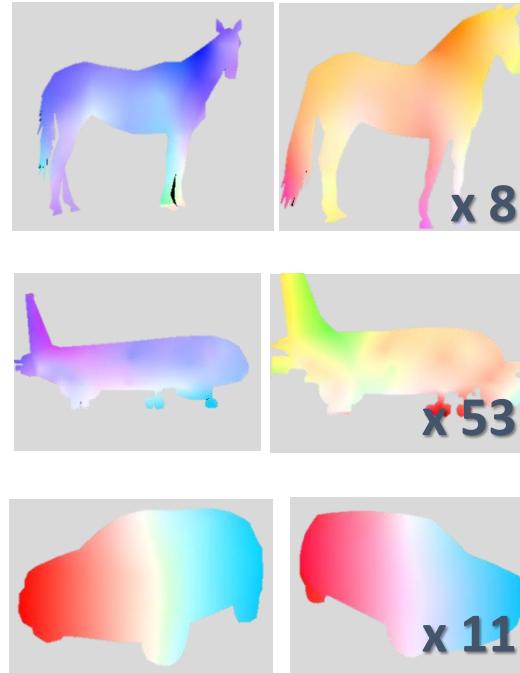
400 pairs

New Dataset

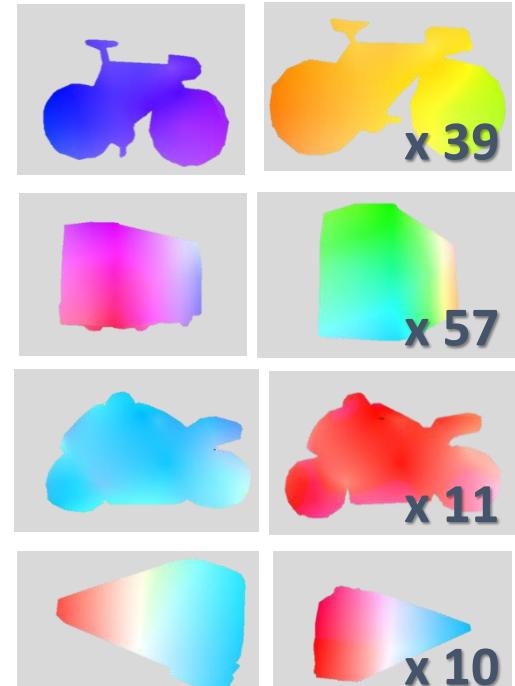
FG3DCar ^[1]



JODS ^[2]



PASCAL ^[3]



Dense ground truth correspondence obtained by interpolating sparse key-point matches (annotated by a user).

400 pairs

Challenges

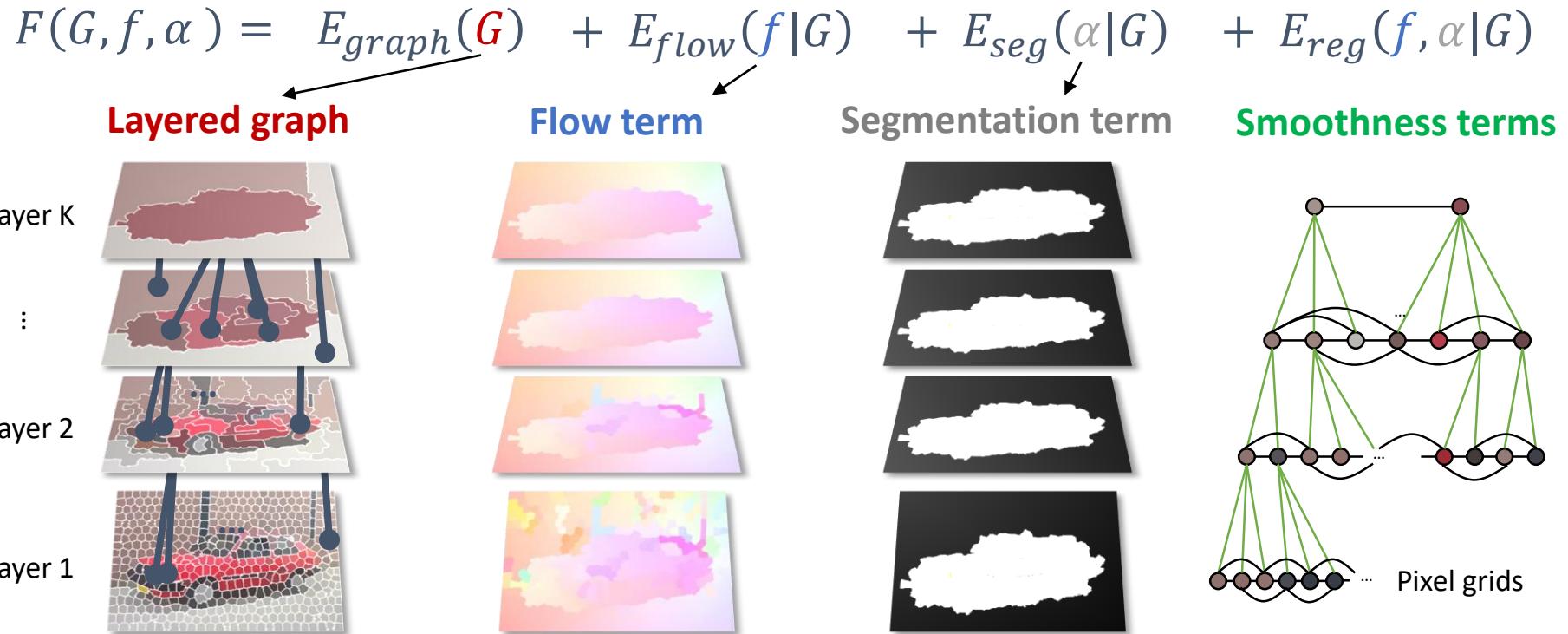
- Objects are unknown
- Appearance, shape similarity cues are weak
- Viewpoints, backgrounds differ

Towards ..
unsupervised visual object discovery

Our Approach

- Jointly recover flow and segmentation
- Hierarchical model
 - (Structure) Layered graph of nested image regions
 - (Continuous Label Space)
 - binary (segmentation)
 - 2D similarity transform (flow) (4-dof)
 - (Spatial regularization)
 - between neighbors
 - between parent-child nodes.
- Energy minimization/Inference
 - Local alpha expansions (graph cuts) [Taniai et al. 2014]

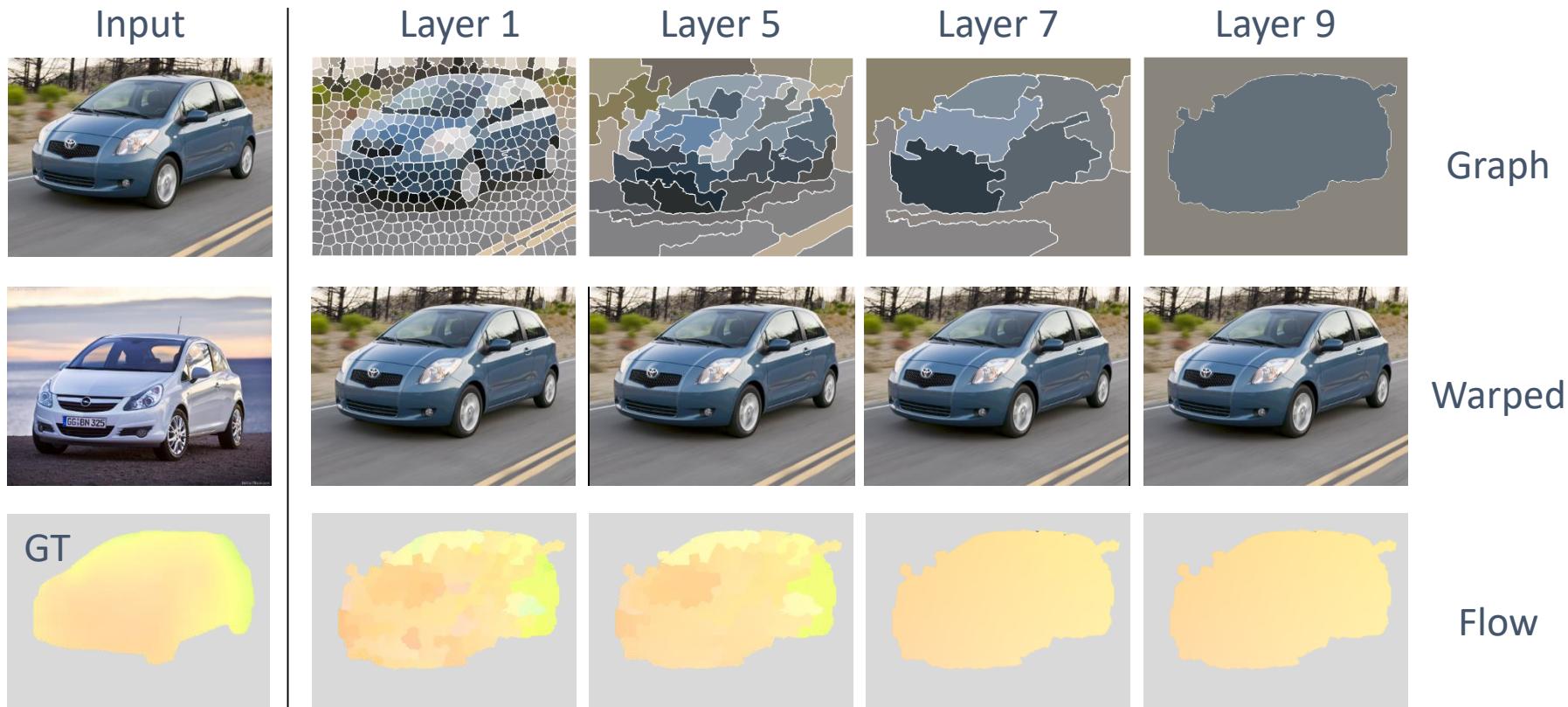
Hierarchical Model



- structure inferred one layer at a time
- Patch matching with HOG descriptors
- FG/BG color likelihoods

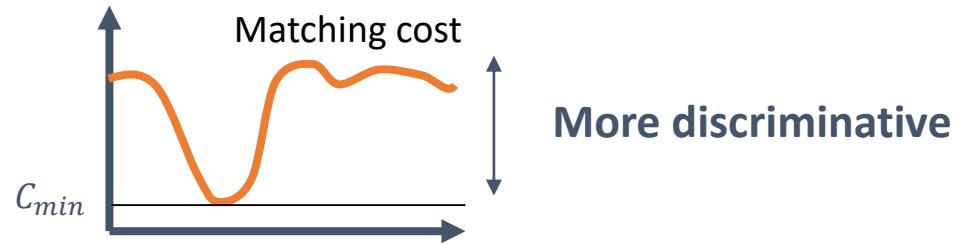
- Spatial neighbor edges
- Parent child edges

Hierarchical flow visualization

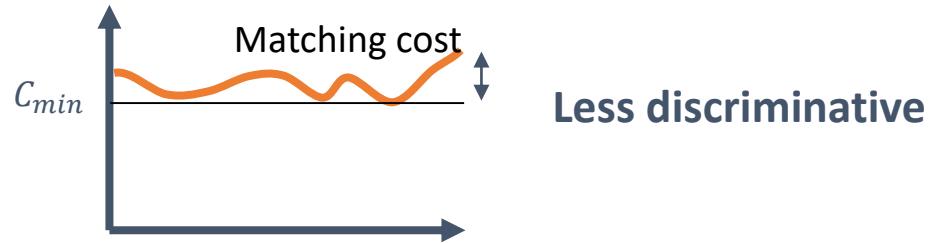


Foreground / Background cues

Foreground patches are likely to have a good match (low cost)

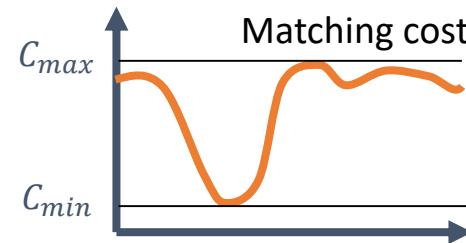
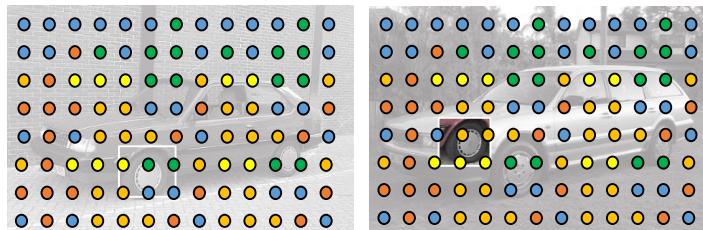


Background patches will have random matches (usually high cost)

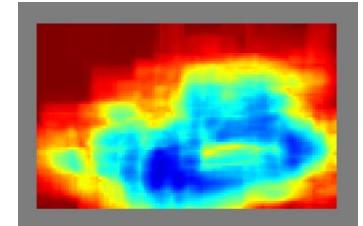


Foreground / Background cues

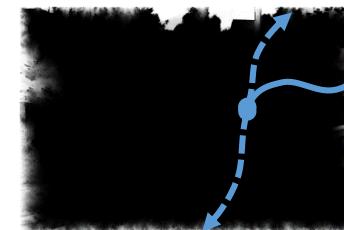
Matching Cost Ratio



$$\text{Ratio} = C_{min}/C_{max}$$



Geodesic distance from image boundary



Shortest path
distance

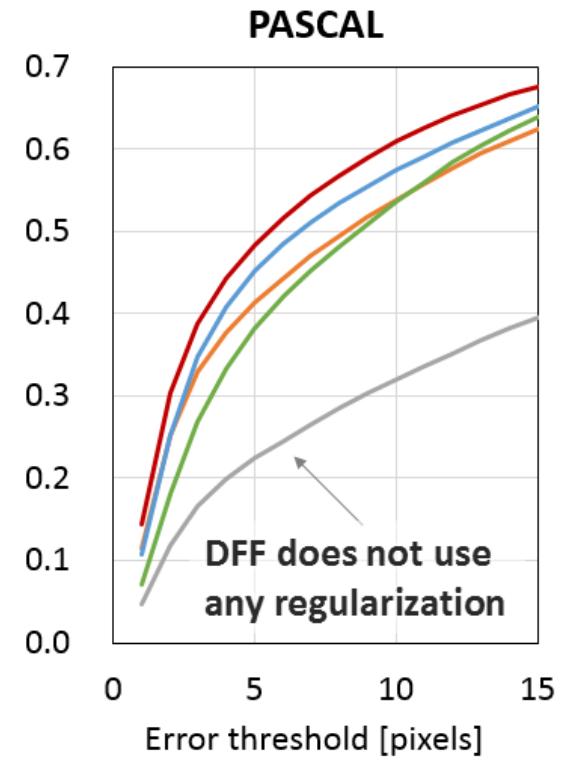
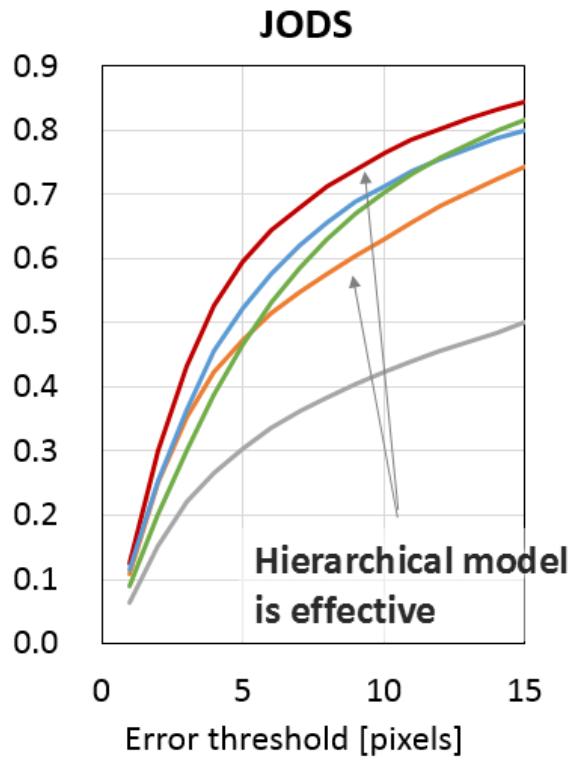
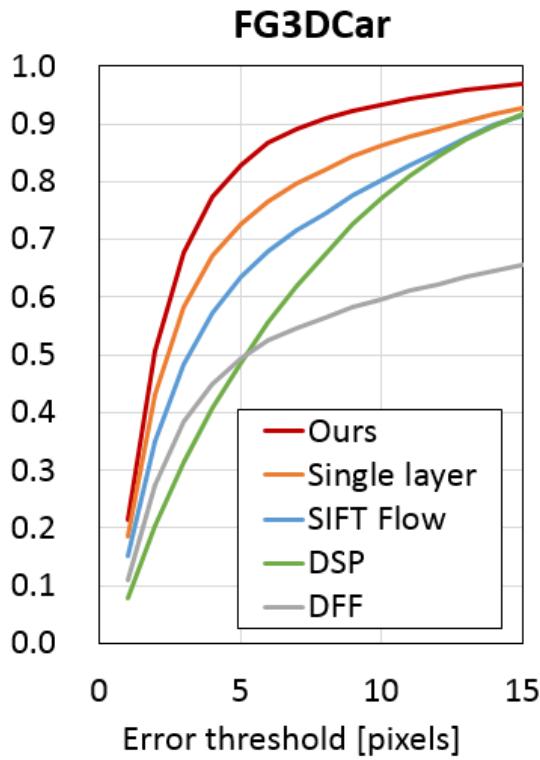
- Construct seeds and initial mask for GrabCut [Rother+04]
- Learn FG/BG color models for each image

Evaluation

Baselines

- Our single layer model
- SIFT Flow [Liu et al. 2008]
- Deformable spatial pyramids (DSP) [Kim et al. 2013]
- DAISY filter flow (DFF) [Yang et al. 2014]
- Cosegmentation by composition [Faktor and Irani 2013]
- Discriminative Clustering [Joulin et al. 2010, 2012]
- NRDC ... [HaCohen et al. 2011]
- Cosegmentation by Co-sketch [Dai et al. 2013]

Flow Accuracy

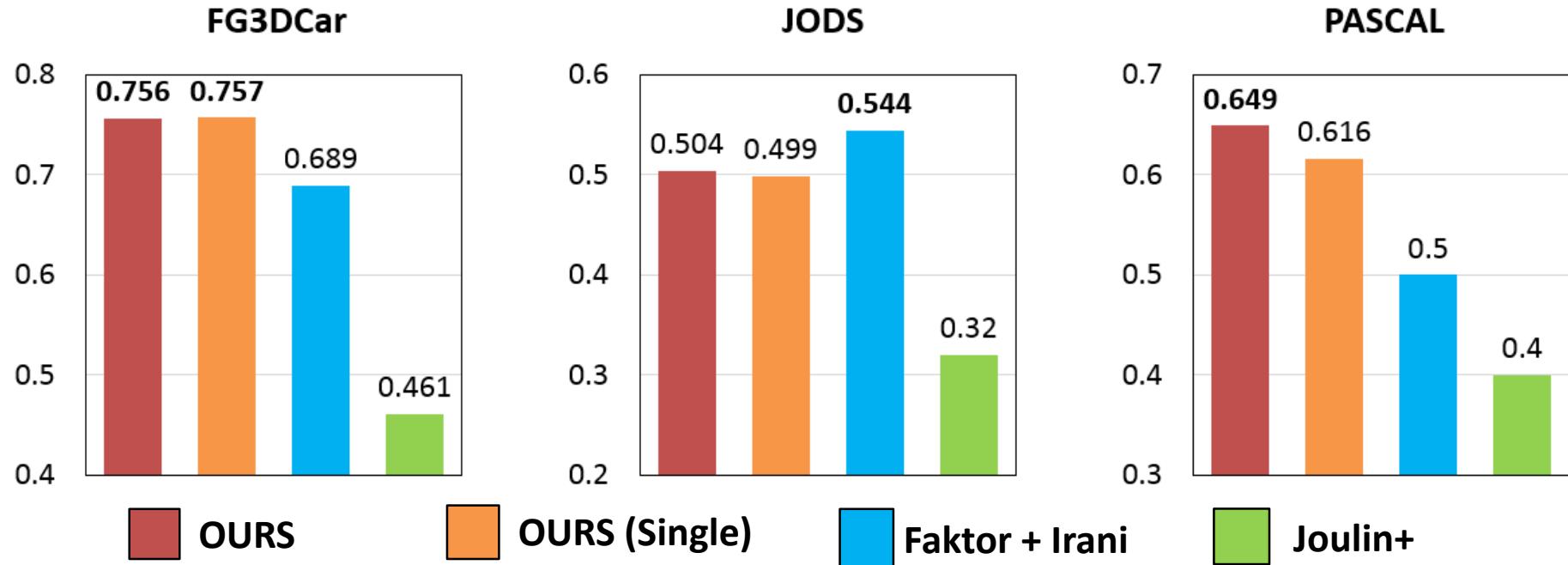


Accuracy Metric

- Percentage of flow errors above a threshold (2d distance)

Our method consistently outperforms all the baselines

Cosegmentation Accuracy

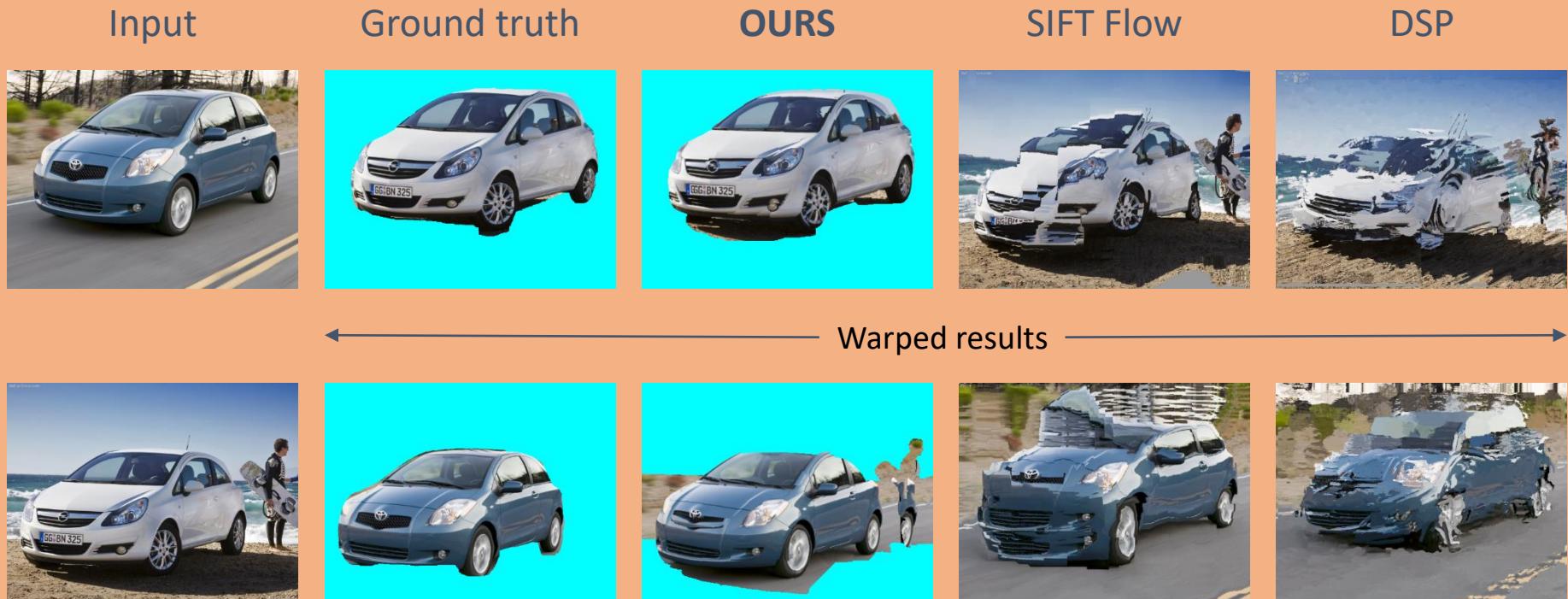


Accuracy Metric

- Intersection-over-union ratio

Our method achieves comparable or better accuracies

Alignment Results



Source: FG3DCar

Alignment Results

Input



Ground truth



OURS



SIFT Flow



DSP



Warped results



Source: JODS

Alignment Results

Input



Ground truth



OURS



SIFT Flow



DSP

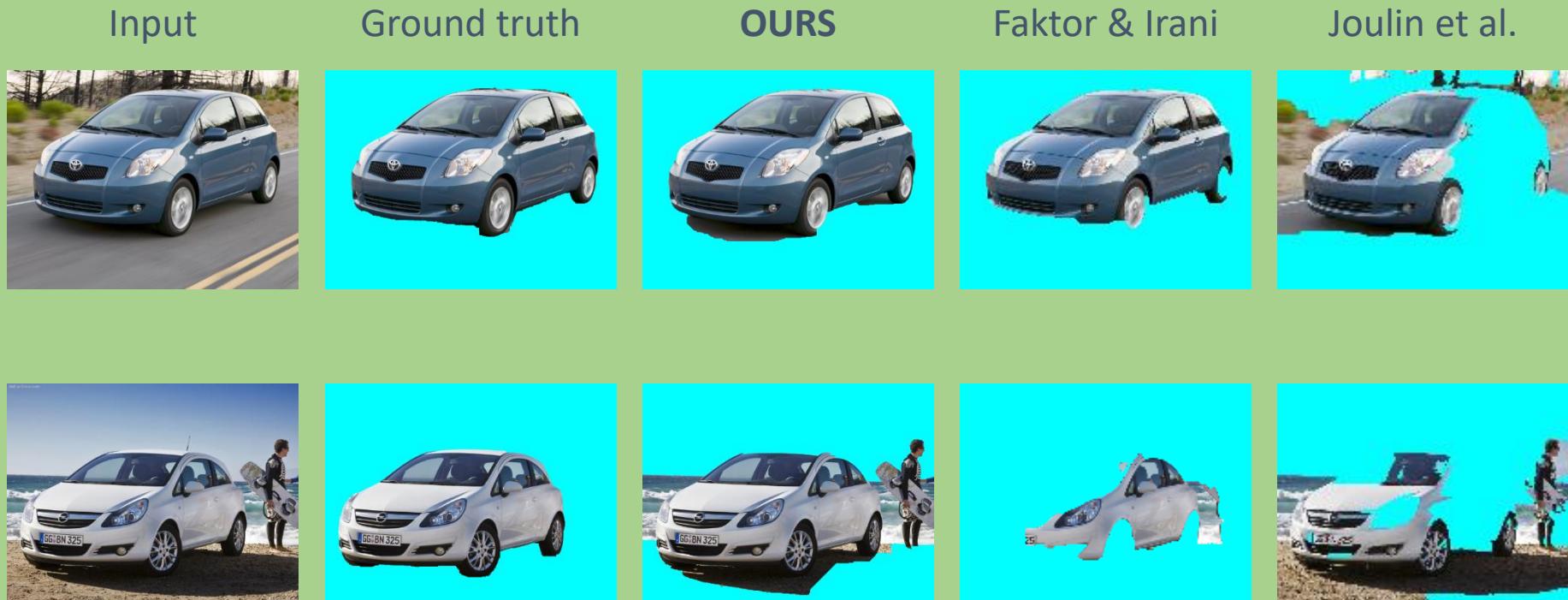


Warped results



Source: PASCAL

Cosegmentation Results



Source: FG3DCar

Cosegmentation Results

Input



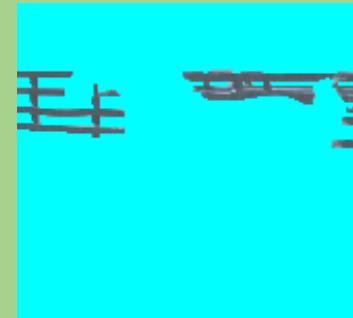
Ground truth



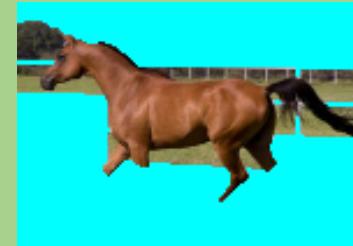
OURS



Faktor & Irani



Joulin et al.



Source: JODS

Cosegmentation Results

Input



Ground truth



OURS



Faktor & Irani



Joulin et al.



Source: PASCAL

Future Work

- Try pre-trained ConvNet features
- Add bi-directional flow consistency
- Use multiple images, add cycle-consistency

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - Color Consistency in Photo Collections

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - **Improved place recognition**
 - Color Consistency in Photo Collections

Leveraging Structure from Motion to Learn Discriminative Codebooks for Scalable Landmark Classification

CVPR 2013



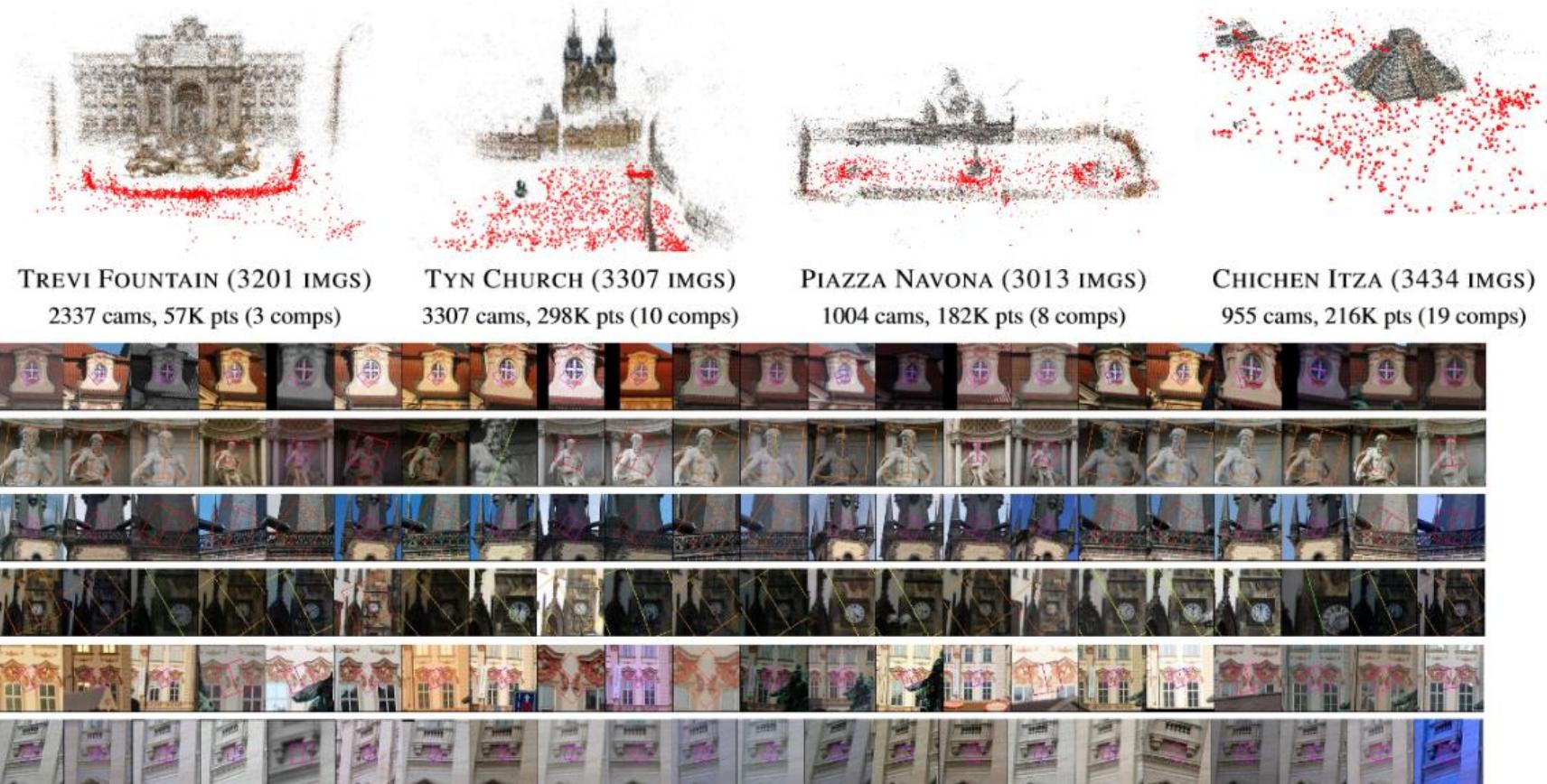
Alessandro Bergamo*
Amazon

* while at Dartmouth College



Lorenzo Torresani
Dartmouth College

SfM from Internet photos



Problem

Goal

A single image from one of k locations. Recognize the location.

Approach

Image categorization (BoW/VLAD/Fisher \rightarrow linear SVM)

Train a binary classifier for each location

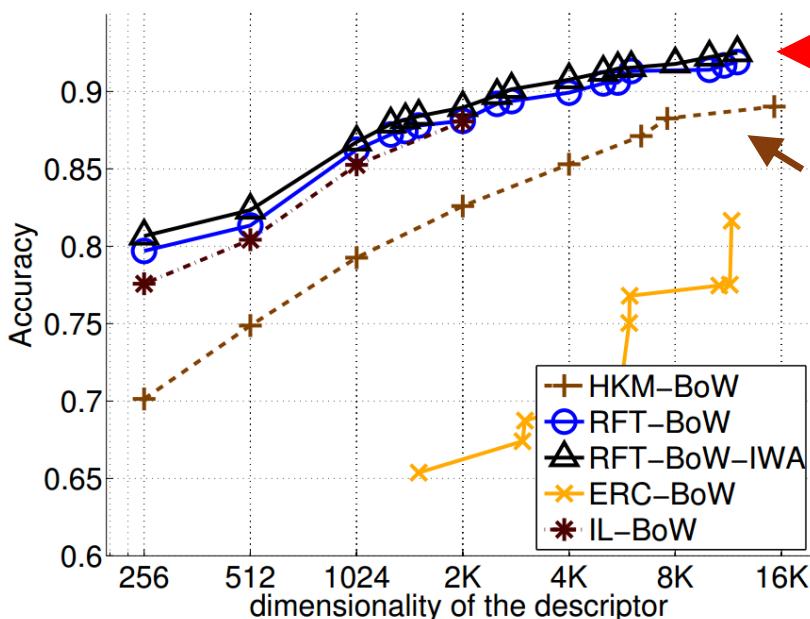
Idea (Discriminative Codebook Learning)

- Each track (n-view correspondence) is a unique class.
- Train a discriminative random forest.
- Use it to quantize/aggregate local descriptors.

Results: top-1 classification accuracy

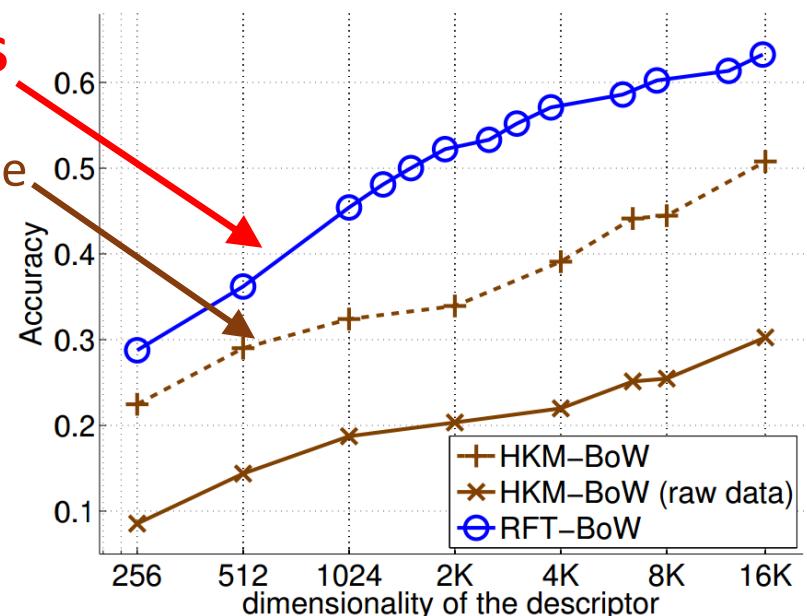
Landmark3D

(25 places, 5K test images)



Landmark-620

(620 places, 62K test images)



- More accurate than k-means (SIFT/DAISY)
- for both BoW/VLAD representations

Overview

- Dense Correspondence Estimation
 - Surface Stereo
 - High Resolution Stereo Matching
- Joint Correspondence and Cosegmentation
 - Align different object instances
- Sparse Correspondences and Applications
 - Improved place recognition
 - **Color Consistency in Photo Collections**

Efficient and Robust Color Consistency for Community Photo Collections

CVPR 2016 (to appear)



Jaesik Park*
Intel Labs



Yu-Wing Tai*
SenseTime

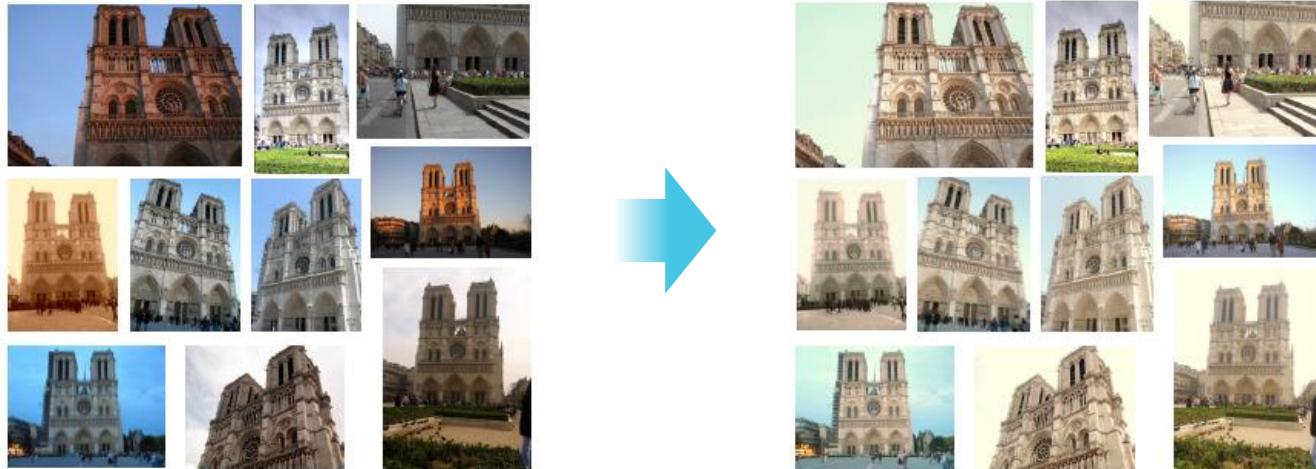


In So Kweon
KAIST

* while at KAIST

Goal

Improve the color consistency of images in a photo collection



Rigid Scenes

Feature matching + Structure from Motion (SfM)

Non-rigid scenes

Feature matching in image pairs

Construct a match graph

Compute maximal cliques [Bron-Kerbosch algorithm (1977)]

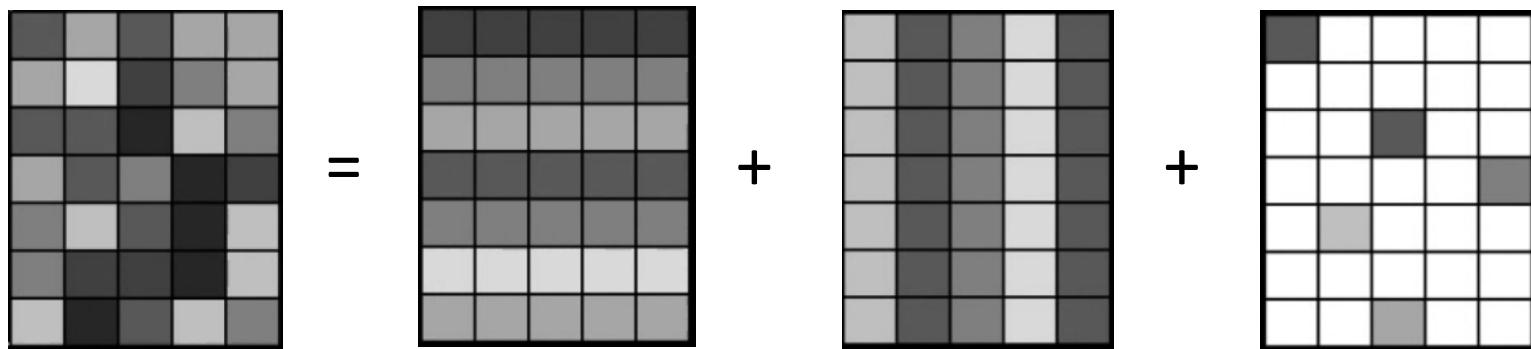
Main Idea

Color Correction Model: $I' = (cI)^\gamma$

Constraints from sparse correspondences ..

$$I_i(x_{ij}) = (c_i a_j e_{ij})^{\gamma_i}$$

Low-rank Matrix Factorization formulation


$$\begin{matrix} \text{Input Matrix} & = & \text{Matrix 1} & + & \text{Matrix 2} & + & \text{Matrix 3} \end{matrix}$$

Low-Rank Matrix Decomposition Technique
(Cabral et al. 2013, in ICCV)

Results – ICE SKATER (36 images)



- Our method is faster than [HaCohen et al. 2012] which requires dense correspondence.
- Robust formulation; resilient to outliers.

Image Stitching



Input Images for Microsoft ICE (stitcher)



Using original images



Using images corrected with
Photoshop CS6



Using our corrected images

Multi-view Stereo

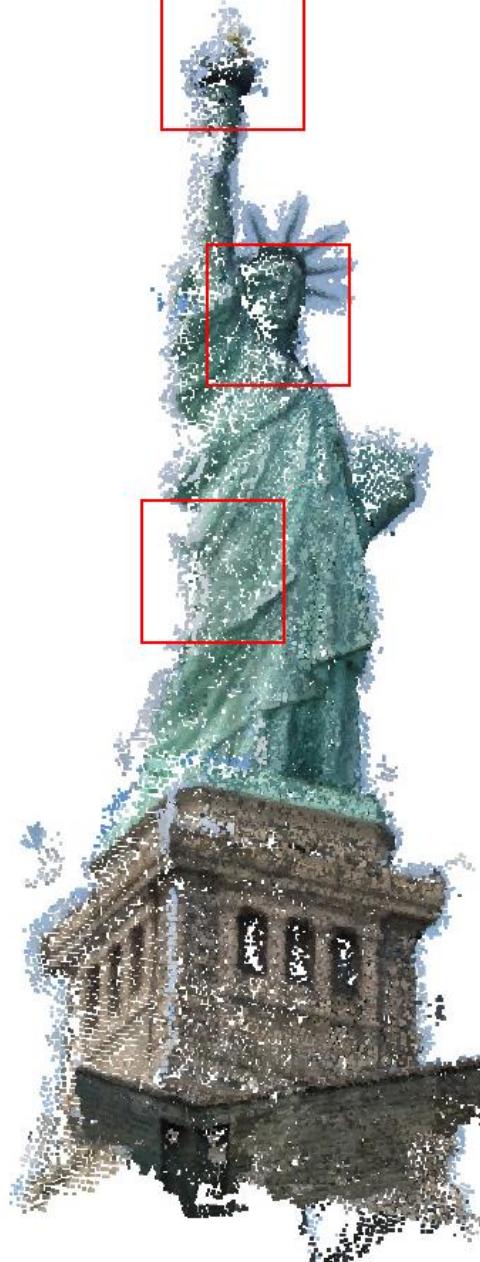
original images



corrected images



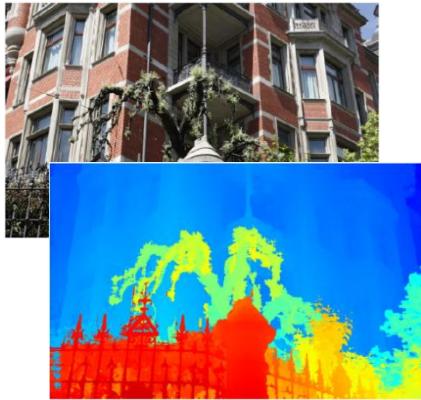
Using original images



Using corrected images



Conclusions



High resolution Stereo

- Local plane sweeps
- Reduce search space
- SGM optimization

Flow + Cosegmentation

- Joint formulation
- Hierarchical MRF model
- Continuous labels
- Graph cuts

Color Consistency in Photo Collections

- Uses sparse feature matches
- Robust matrix factorization
- Efficient color transfer