

Looking Back, Moving Forward



Computational Photography
Derek Hoiem, University of Illinois

Today

- Requested topics
 - 3D reconstruction
 - Light transport
 - Event cameras
- Beyond this class...
- ICES forms
- Reminder: final project
 - Write-ups due Dec 15 11:59pm
 - Poster presentations on Dec 16 at 8am on the first floor of Siebel
 - Half of class will present at one time, then switch
 - Everyone is assigned to review two posters (and should also look at the others that are of interest)
 - See Piazza pinned post for updates

Project 5

- Incomplete list of excellent projects

<https://cyu17.web.illinois.edu/cs445/proj5/> - many bells and whistles, inserted mickey mouse

<https://aipark2.web.illinois.edu/cs445/proj5/> - nice page, student union video

<https://lehan2.web.illinois.edu/cs445/proj5/> - add neon lights

<https://aayushr2.web.illinois.edu/cs445/proj5/> - remove camera shake, bus video

<https://susiel2.web.illinois.edu/cs445/proj5/> - seam-based blending

<https://yicheng9.web.illinois.edu/cs445/proj5/> - added duck, reduced camera shake

<https://chan104.web.illinois.edu/cs445/proj5/> - blending, reduced camera shake

This course has provided fundamentals

- How photographs are captured from and relate to the 3D scene
- How to think of an image as: a signal to be processed, a graph to be searched, an equation to be solved
- How to manipulate photographs: cutting, growing, compositing, morphing, stitching
- Basic principles of computer vision: filtering, correspondence, alignment

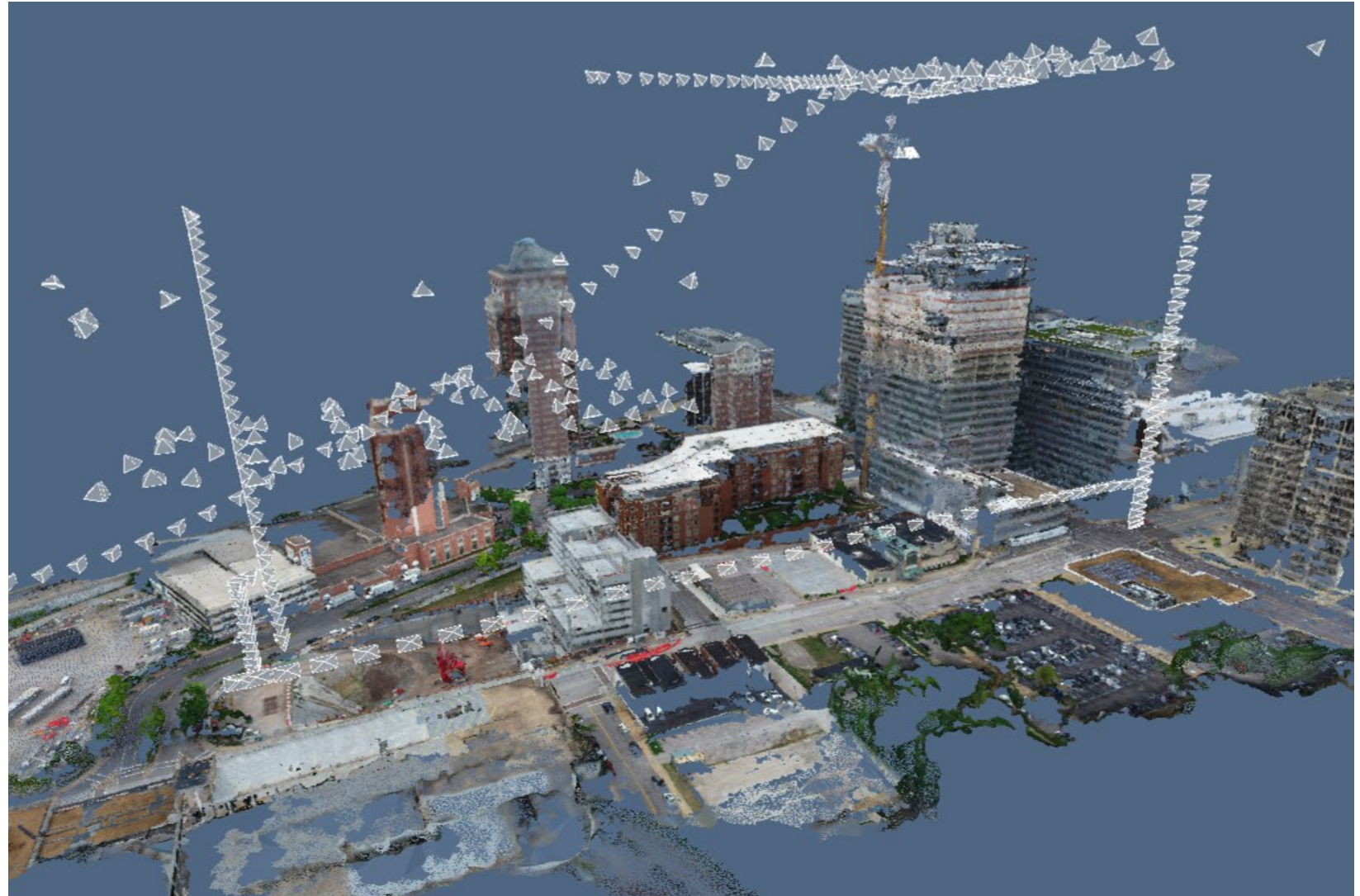
What else is out there?

Lots!

- Machine learning
- Videos and motion
- 3D reconstruction
- Scene understanding
- Better/cheaper devices
- ...

How to create 3D model from multiple images

1. Solve for camera poses
2. Propose and verify 3D points by matching
3. Fit a surface to the points



Incremental Structure from Motion (SfM)

Goal: Solve for camera poses and 3D points in scene



Incremental SfM

1. Compute features

2. Match images

3. Reconstruct

a) Solve for pose and 3D points in two cameras

b) Solve for pose of additional camera(s) that observe reconstructed 3D points

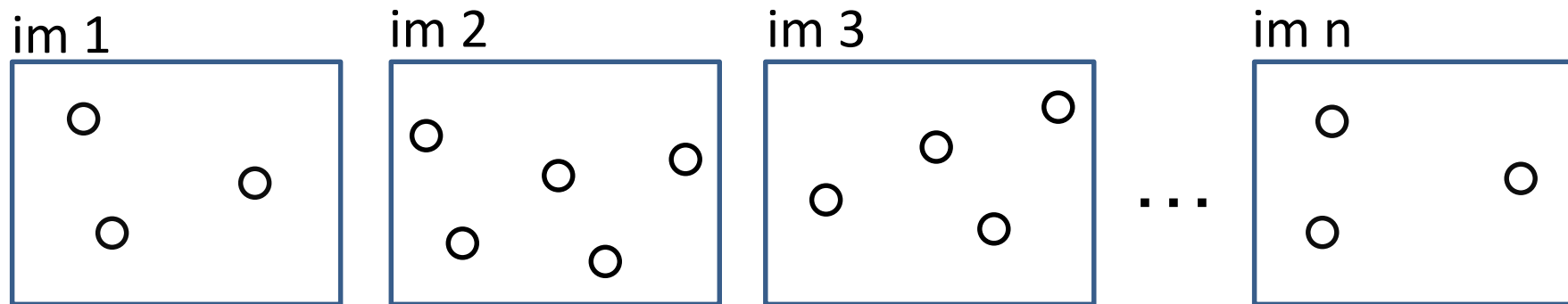
c) Solve for new 3D points that are viewed in at least two cameras

d) Bundle adjust to minimize reprojection error



Incremental SFM: **detect features**

- Feature types: SIFT, ORB, Hessian-Laplacian, ...



Each circle represents a set of detected features

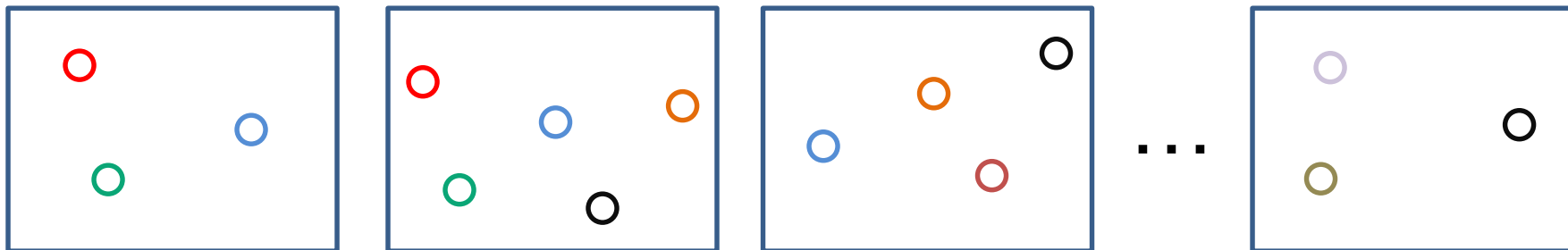
Incremental SFM: match features and images

For each pair of images:

1. Match feature descriptors via approximate nearest neighbor and apply Lowe's ratio test
2. Solve for F and find inlier feature correspondences

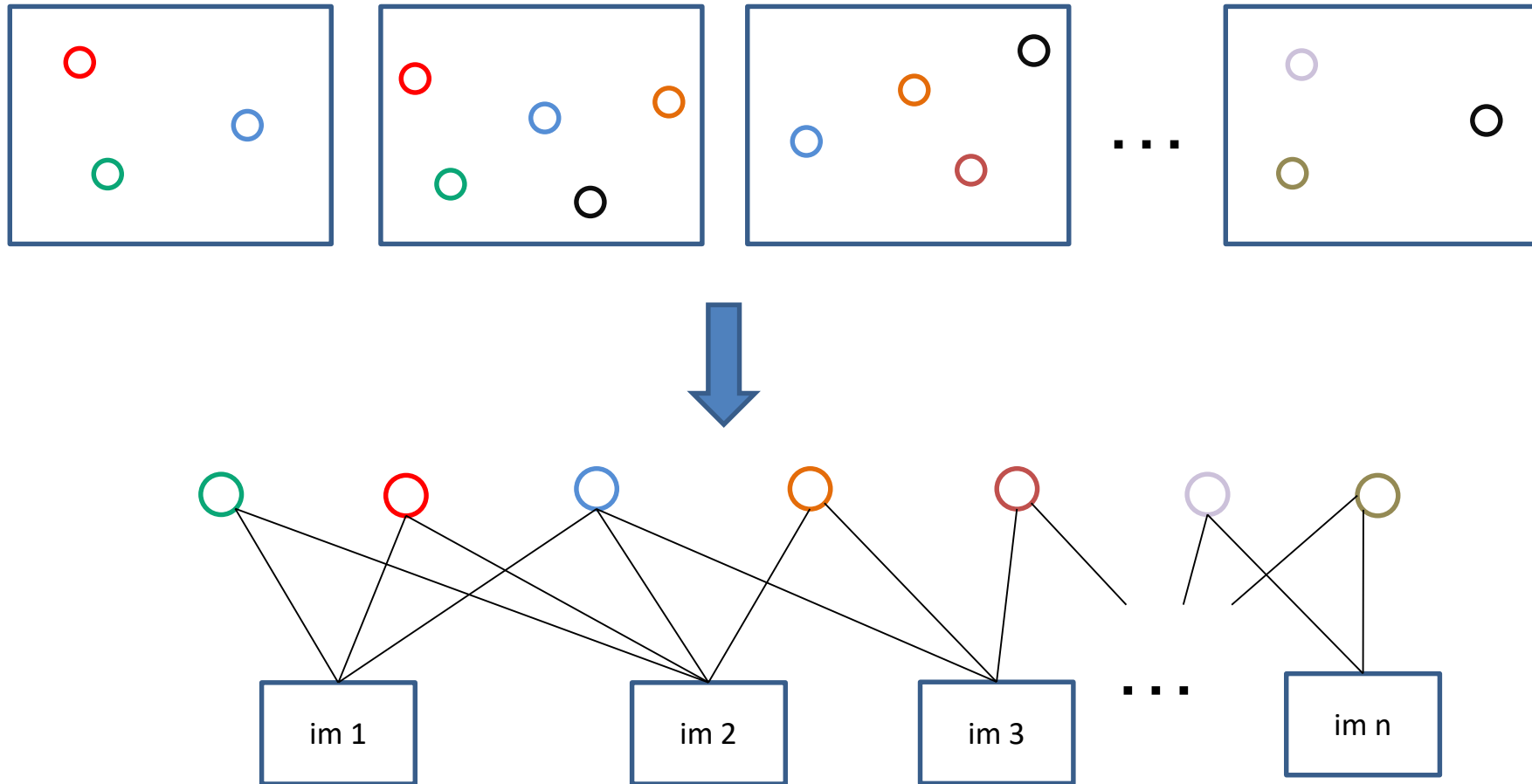
- Speed tricks

- Use vocabulary tree to get image match candidates
- Use GPS coordinates to get match candidates, if available



Points of same color have been matched to each other

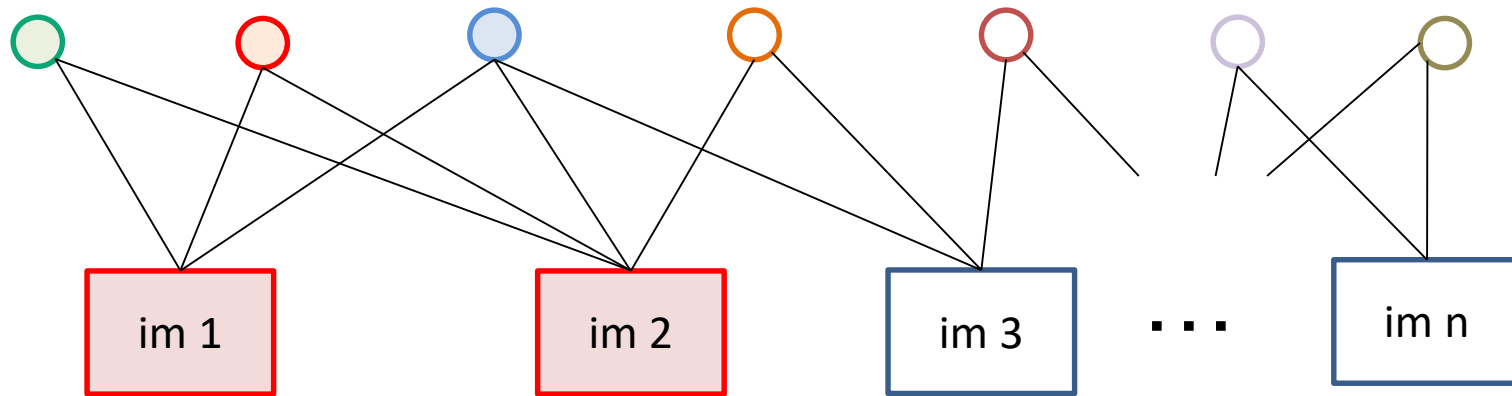
Incremental SFM: create tracks graph



tracks graph: bipartite graph between observed 3D points and images

Incremental SFM: initialize reconstruction

1. Choose two images that are likely to provide a stable estimate of relative pose
 - E.g., $\frac{\# \text{ inliers for } H}{\# \text{ inliers for } F} < 0.7$ and many inliers for F
2. Get focal lengths from EXIF, estimate essential matrix using 5-point algorithm, extract pose R_2, t_2 with $R_1 = \mathbf{I}, t_1 = \mathbf{0}$
3. Solve for 3D points given poses
4. Perform bundle adjustment to refine points and poses

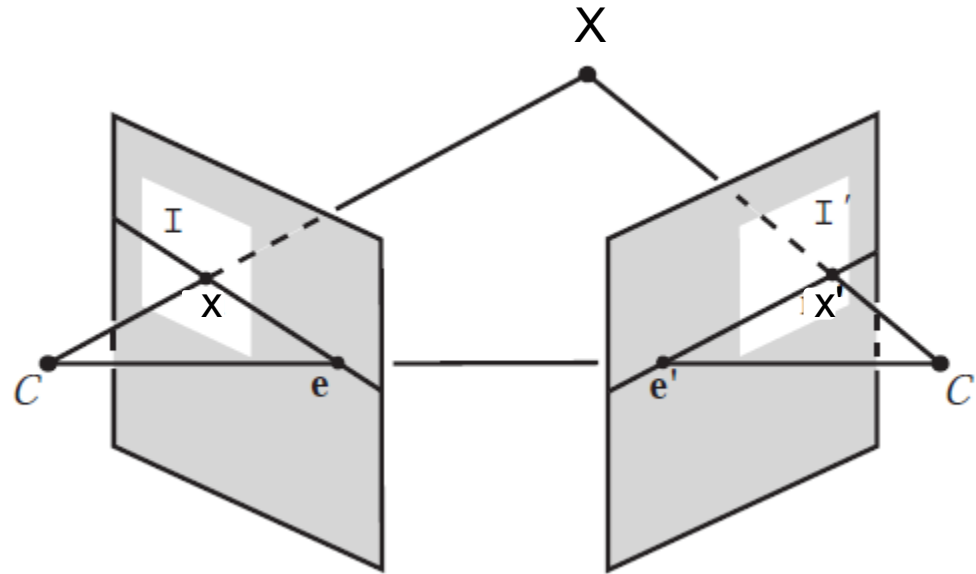


filled circles = “triangulated” points

filled rectangles = “resectioned” images (solved pose)

Triangulation: Linear Solution

- Generally, rays $C \rightarrow x$ and $C' \rightarrow x'$ will not exactly intersect
- Can solve via SVD, finding a least squares solution to a system of equations



$$\mathbf{x} = \mathbf{P}\mathbf{X}$$

$$\mathbf{x}' = \mathbf{P}'\mathbf{X}$$



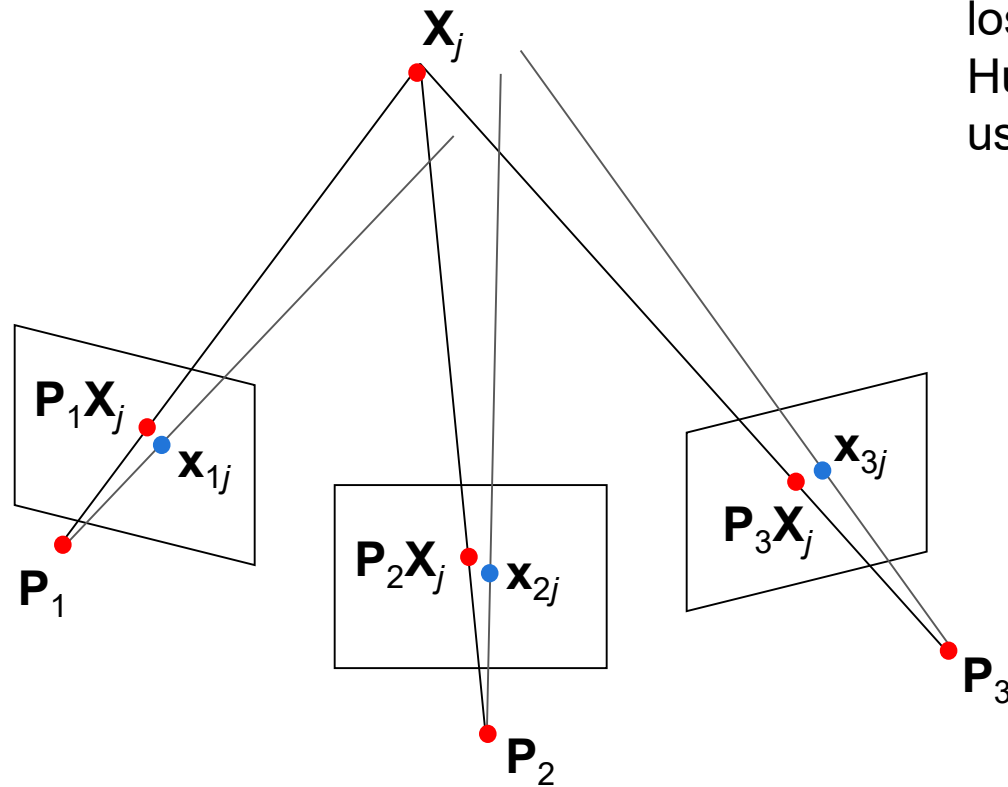
$$\mathbf{A}\mathbf{X} = \mathbf{0} \quad \mathbf{A} = \begin{bmatrix} u\mathbf{p}_3^T - \mathbf{p}_1^T \\ v\mathbf{p}_3^T - \mathbf{p}_2^T \\ u'\mathbf{p}'_3{}^T - \mathbf{p}'_1{}^T \\ v'\mathbf{p}'_3{}^T - \mathbf{p}'_2{}^T \end{bmatrix}$$

Bundle adjustment

- Non-linear method for refining structure and motion
- Minimizing reprojection error

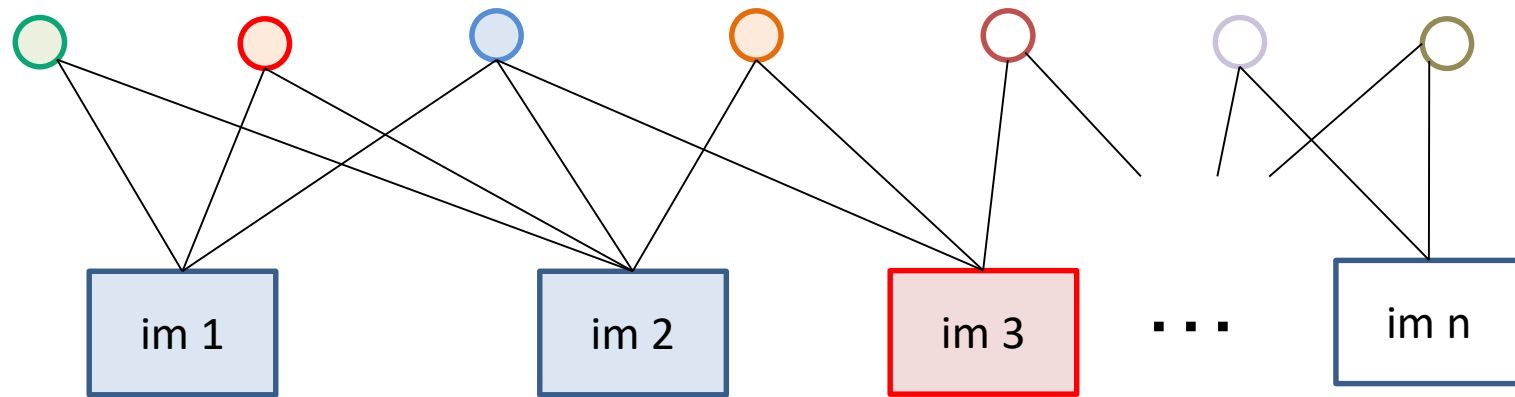
$$E(\mathbf{P}, \mathbf{X}) = \sum_{i=1}^m \sum_{j=1}^n D(\mathbf{x}_{ij}, \mathbf{P}_i \mathbf{X}_j)^2$$

Often a robust loss, such as Huber loss is used



Incremental SFM: grow reconstruction

1. Resection: solve pose for image(s) that have the most triangulated points
2. Triangulate: solve for any new points that have at least two cameras
3. Remove 3D points that are outliers
4. Bundle adjust
 - For speed, only do full bundle adjust after some percent of new images are resectioned
5. Optionally, align with GPS from EXIF or ground control points (GCP)

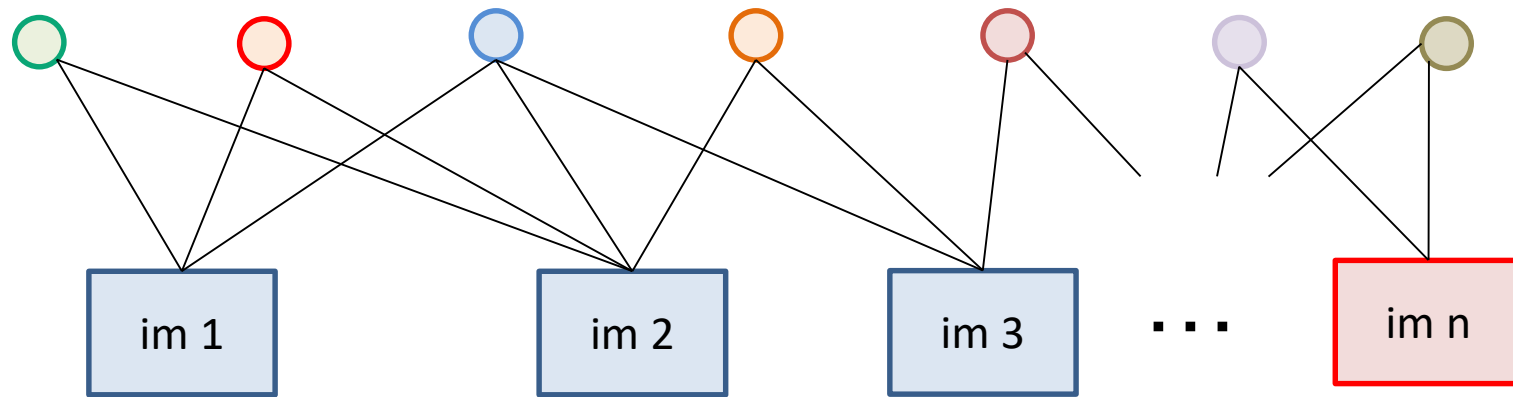


filled circles = “triangulated” points

filled rectangles = “resectioned” images (solved pose)

Incremental SFM: grow reconstruction

1. Resection: solve pose for image(s) that have the most triangulated points
2. Triangulate: solve for any new points that have at least two cameras
3. Remove 3D points that are outliers
4. Bundle adjust
 - For speed, only do full bundle adjust after some percent of new images are resectioned
5. Optionally, align with GPS from EXIF or ground control points (GCP)



filled circles = “triangulated” points

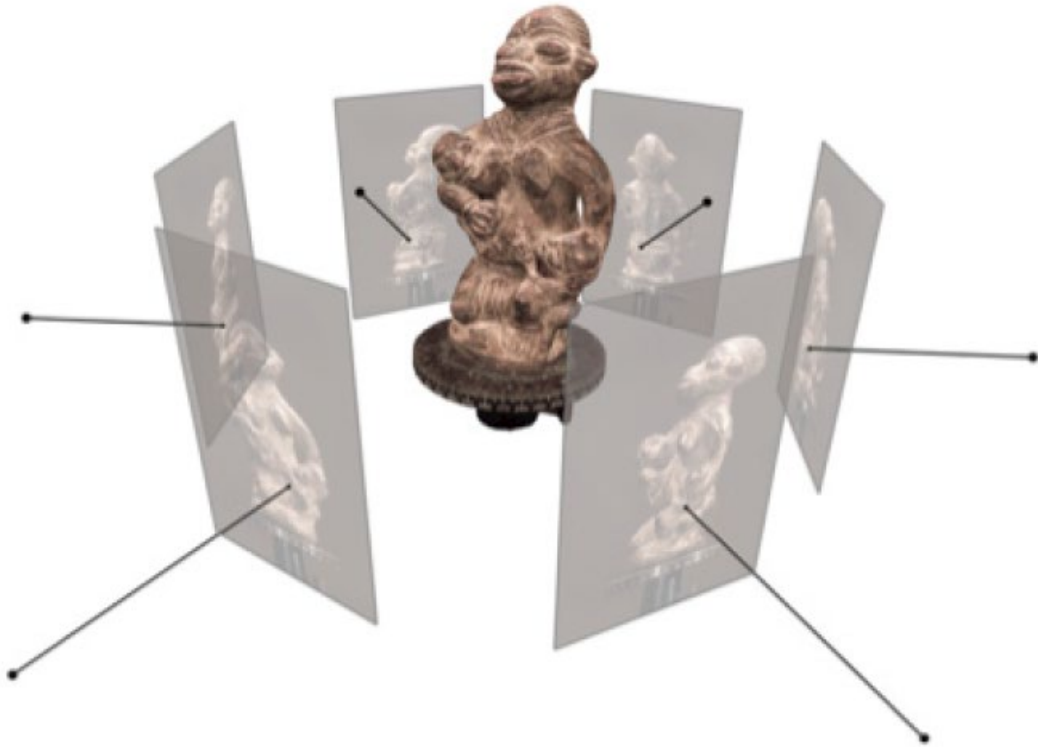
filled rectangles = “resectioned” images (solved pose)

Important recent papers and methods for SfM

- Snaveley thesis (2008): intro to SfM in Chapter 3
- Visual SfM: Visual SfM (Wu 2013)
 - Used to be the best incremental SfM software (but not anymore and closed source); paper still very good
- COLMAP
 - Good open source system based on “Structure-from-motion revisited” (Schonberger Frahm 2016)
- OpenSfM:
 - Python open-source system, easy to read and modify

Reconstruction of Cornell (Crandall et al. ECCV 2011)

Multiview Stereo: propose and verify 3D points by matching pixel patches across images



Select depth at each pixel that minimizes NCC of patches with other images

Key Assumptions

- Enough texture to match
- Surface looks the same from each view (non-reflective)

Multiview Stereo: recommended reading

“Multiview Stereo: a tutorial” by Yasu Furukawa

http://www.cse.wustl.edu/~furukawa/papers/fnt_mvs.pdf

COLMAP:

- Code based on “Pixelwise View Selection for Unstructured Multi-View Stereo” by Schonberger et al. 2016

Surface Reconstruction

Floating scale surface reconstruction:

<http://www.gcc.tu-darmstadt.de/home/proj/fssr/>

Constrained Delaunay triangulation

- Create 3D triangulation of dense points and remove faces that conflict with observed points

Deep Image Prior

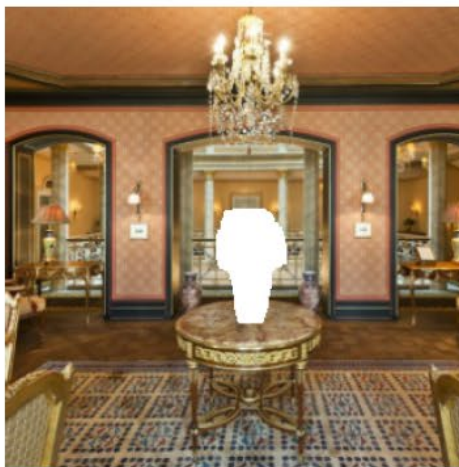
Dmitry Ulyanov
Skolkovo Institute of Science
and Technology, Yandex
dmitry.ulyanov@skoltech.ru

Andrea Vedaldi
University of Oxford
vedaldi@robots.ox.ac.uk

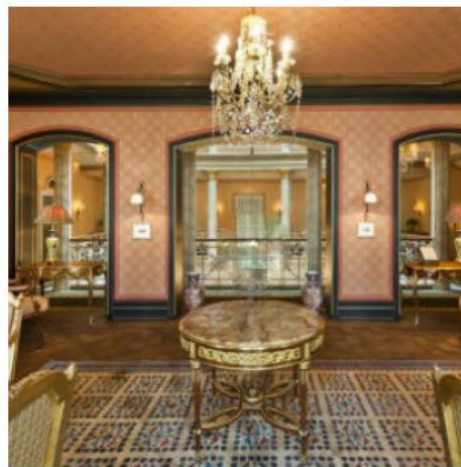
Victor Lempitsky
Skolkovo Institute of Science
and Technology (Skoltech)
lempitsky@skoltech.ru

Surprising result: A randomly initialized decoder network, when trained to reproduce a corrupted image, fixes the noise, holes, etc.

The network structure acts as a prior!



(a) Corrupted image

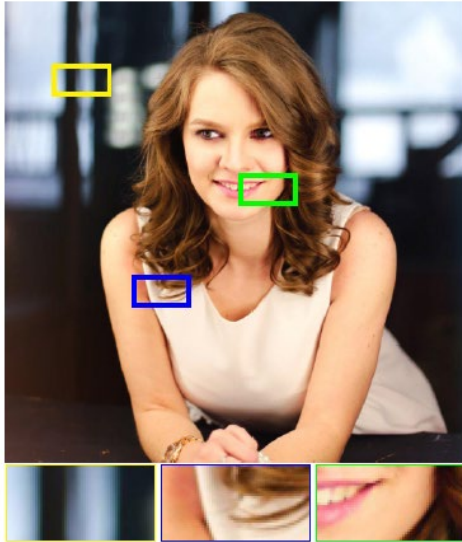


(b) Global-Local GAN [15]

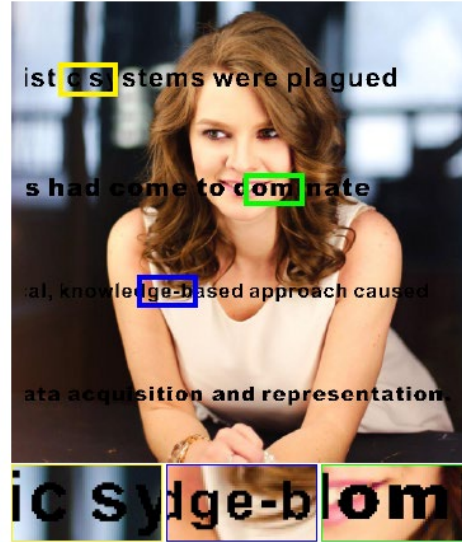


(c) Ours, LR = 0.01

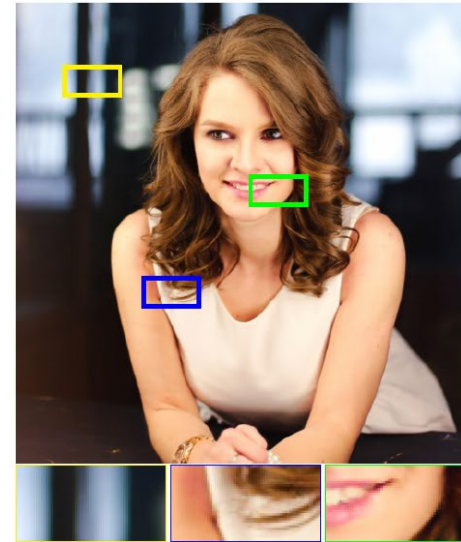
Magic or math? Gradient descent on encoder network to reproduce Original produces a cleaner image. Even better than recent methods designed to solve this problem.



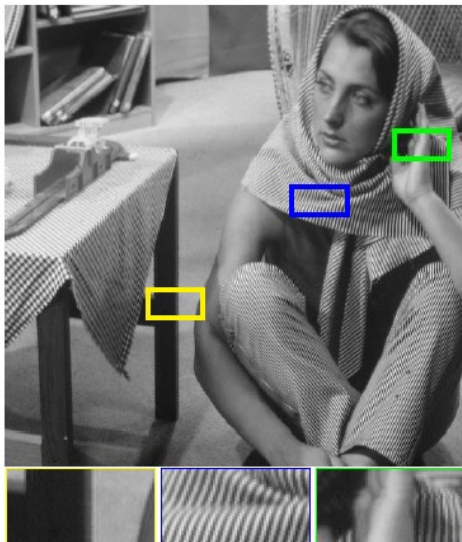
(a) Original image



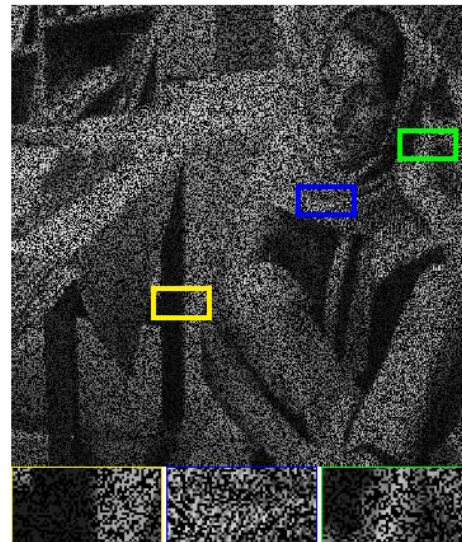
(b) Corrupted image



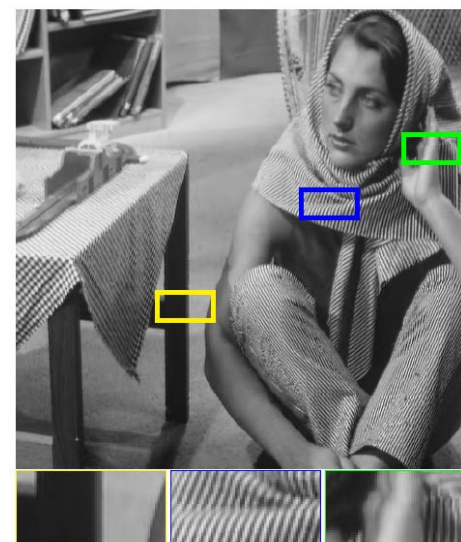
(d) Deep Image Prior



(e) Original image



(f) Corrupted image



(h) Deep Img. Prior, PSNR = 32.22

Computational Mirrors: Blind Inverse Light Transport by Deep Matrix Factorization

NIPS 2019

Miika Aittala
MIT
miika@csail.mit.edu

Prafull Sharma
MIT
prafull@mit.edu

Lukas Murmann
MIT
lmurmann@mit.edu

Adam B. Yedidia
MIT
adamy@mit.edu

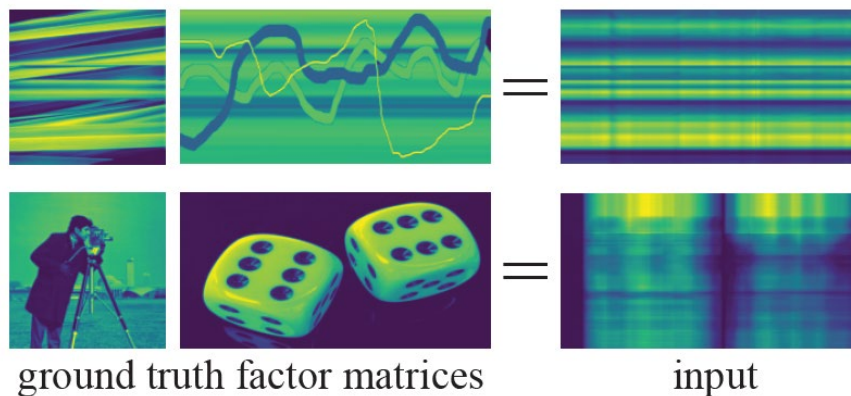
Gregory W. Wornell
MIT
gww@mit.edu

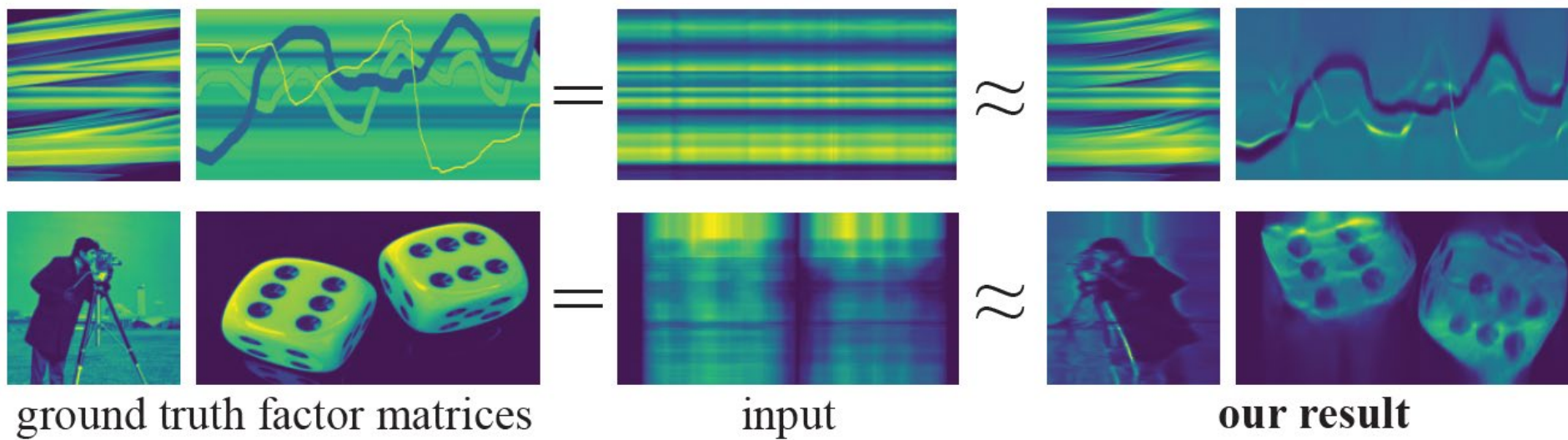
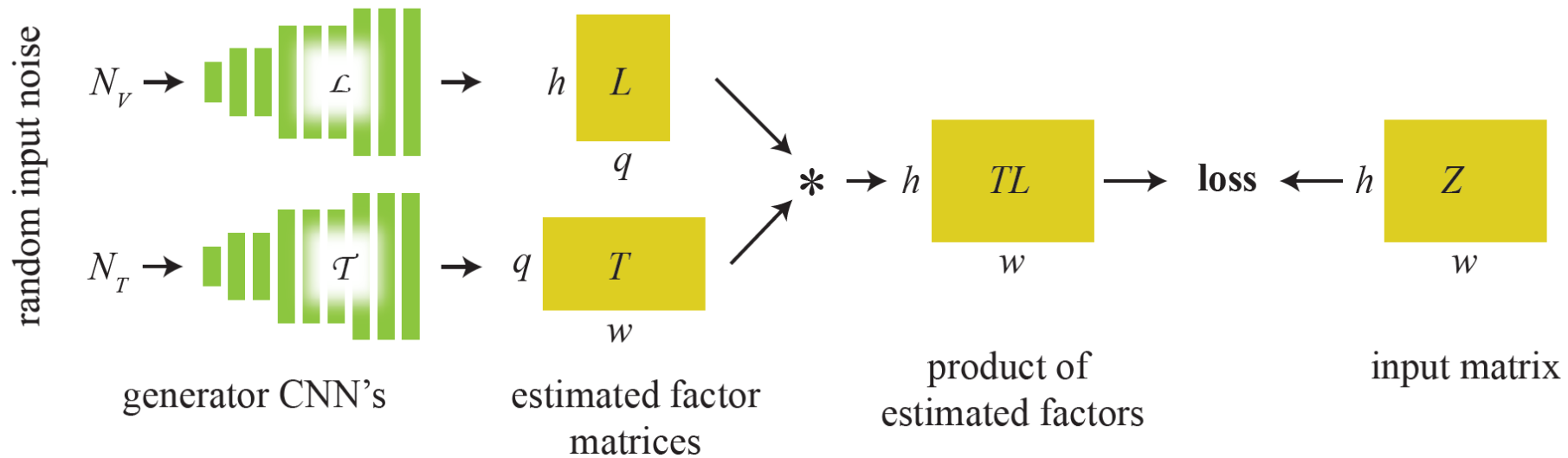
William T. Freeman
MIT, Google Research
billf@mit.edu

Frédo Durand
MIT
fredo@mit.edu

Now take it a step further. If you have the product of two images, you can recover the factors.

Note: there are practically infinitely many useless solutions to this problem.







- Each “pixel” of light on the projector lights the scene, producing an image
- The total image is the sum of images from each pixel.
- Observed image can be factorized into surface colors and projected image (assuming no ambient light)

<https://www.youtube.com/watch?v=bzsfREU2dDM>

Event cameras

- First commercially produced in 2008
- Respond only when individual pixels change intensity
 - Corresponds to camera or scene motion
- 1 micro-second latency
- High dynamic range
- 100x less power than standard camera

Overview: <https://www.youtube.com/watch?v=LauQ6LWTkxM>

3D Reconstruction: <https://www.youtube.com/watch?v=fA4MiSzYHWA>

Two more

- Handwriting beautification (Zitnick SG'13)
 - Example of user assistance
- Semantic image synthesis (Park et al. CVPR 2019)

Trends and Future of Computational Photography

- Camera phones continue to serve as a platform for latest advances in hardware and software
 - Depth may be commonly available
- VR / AR blend graphics with tracking and understanding of environment
 - Killer app outside of games and teleconferencing?
- Design smart programs that work together with people
 - This is #1 from Harry Shum, Exec VP of AI and Research at Microsoft

How can you learn more?

- Relevant courses
 - Production graphics (CS 419)
 - Machine learning (CS 446 and others)
 - Deep learning
 - Computer vision (CS 543)
 - Optimization methods (CS 544)
 - Parallel processing / GPU
 - HCI, data mining, NLP, robotics

How can you learn more?

- Conference proceedings
 - Vision: CVPR, ICCV, ECCV, NIPS
 - Computational photography: ICCP
 - Graphics: SIGGRAPH, SIGGRAPH Asia

Computer Vision (with Prof Gupta Spring 2020)

Similar stuff to CP

- Camera models, filtering, single-view geometry, light and capture

New stuff

- Mid-level vision
 - Edge detection, clustering, segmentation
- Machine learning
- Recognition
 - Image features and classifiers
 - Object category recognition
 - Action/activity recognition
- Videos
 - Tracking, optical flow
 - Structure from motion
- Multi-view geometry

How do you learn more?

Explore and fiddle!

Thank you!

ICES forms

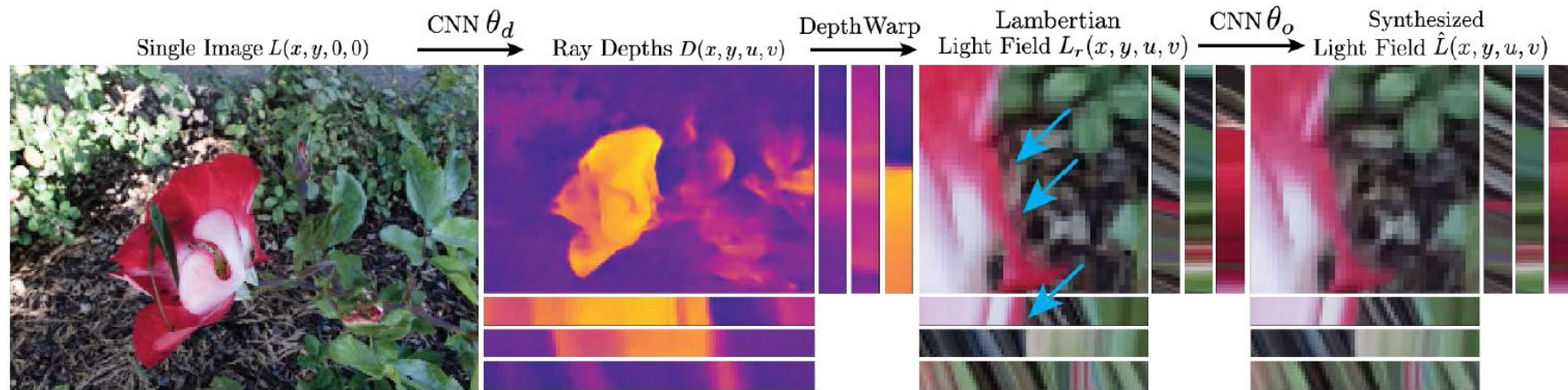
Image \rightarrow Light Field

Learning to Synthesize a 4D RGBD Light Field from a Single Image

Pratul P. Srinivasan¹, Tongzhou Wang¹, Ashwin Sreelal¹, Ravi Ramamoorthi², Ren Ng¹

¹University of California, Berkeley

²University of California, San Diego



<https://www.youtube.com/watch?v=yLCvWoQLnms>

Superresolution

EnhanceNet: Single Image Super-Resolution Through Automated Texture Synthesis

Mehdi S. M. Sajjadi Bernhard Schölkopf Michael Hirsch



Bicubic

ENet-E

ENet-PAT

Ground Truth

E: Optimize least squares objective with upsampling network

PAT: Optimize “perceptual” (VGG features) loss, adversarial loss, texture corr loss



(a) Input

(b) SR [18]

(c) SR [18]+Deblur [33]

(d) Deblur [33]

(e) Deblur [33]+SR [18]

(f) Ours

(g) GT

Learning to Super-Resolve Blurry Face and Text Images

Pretty similar to above, more limited domain

Xiangyu Xu^{1,2,3} Deqing Sun^{3,4} Jinshan Pan⁵ Yujin Zhang¹

Hanspeter Pfister³ Ming-Hsuan Yang²

¹Tsinghua University ²University of California, Merced ³Harvard University

⁴Nvidia ⁵Nanjing University of Science & Technology

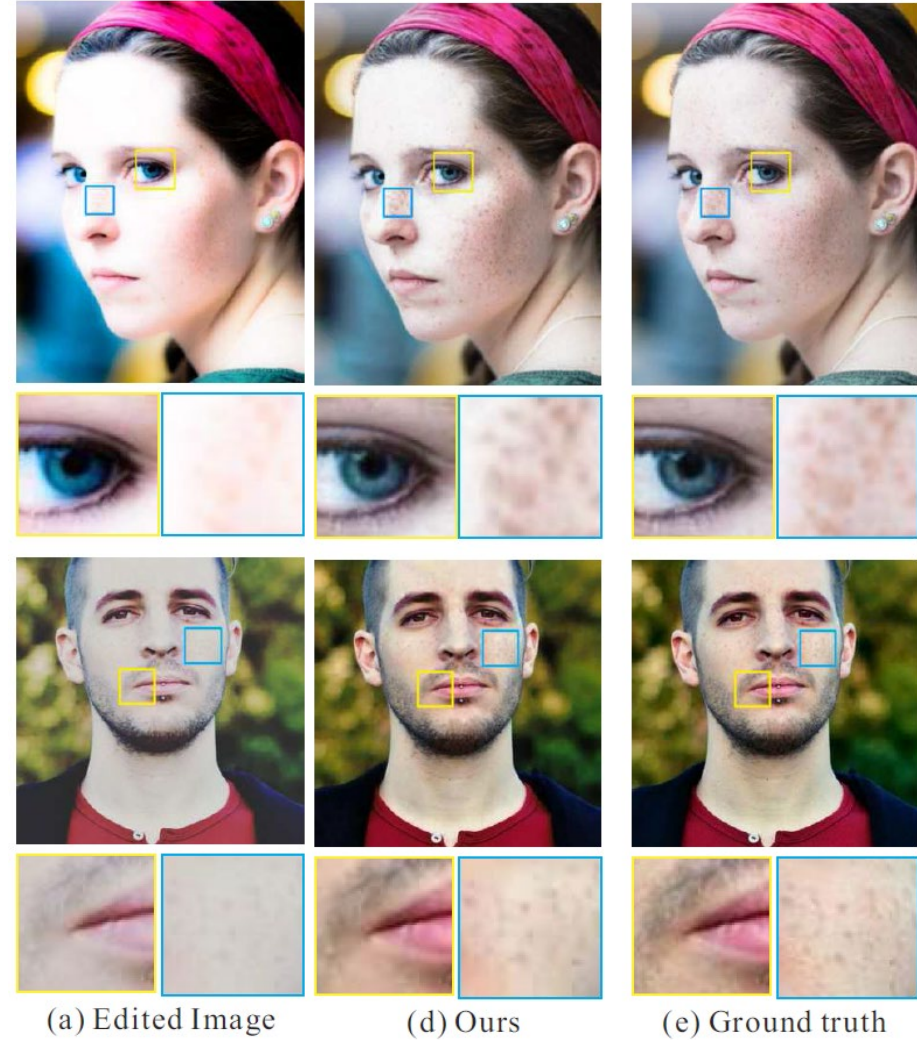
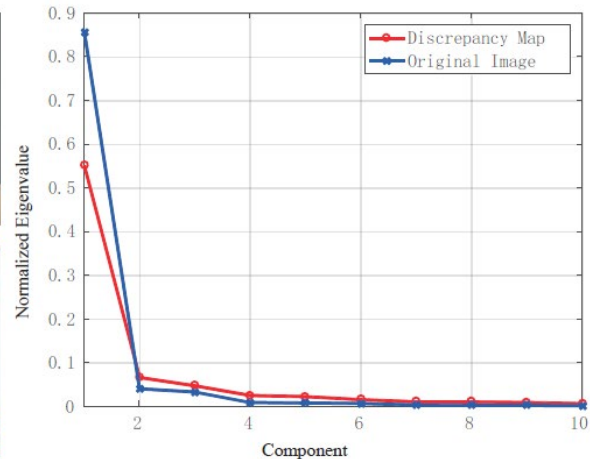
De-beautification

Makeup-Go: Blind Reversion of Portrait Edit*

Ying-Cong Chen¹ Xiaoyong Shen² Jiaya Jia^{1,2}

¹The Chinese University of Hong Kong ²Tencent Youtu Lab

ycchen@cse.cuhk.edu.hk dylanshen@tencent.com leojia9@gmail.com



Network regresses principal components of discrepancy map

LDR --> HDR

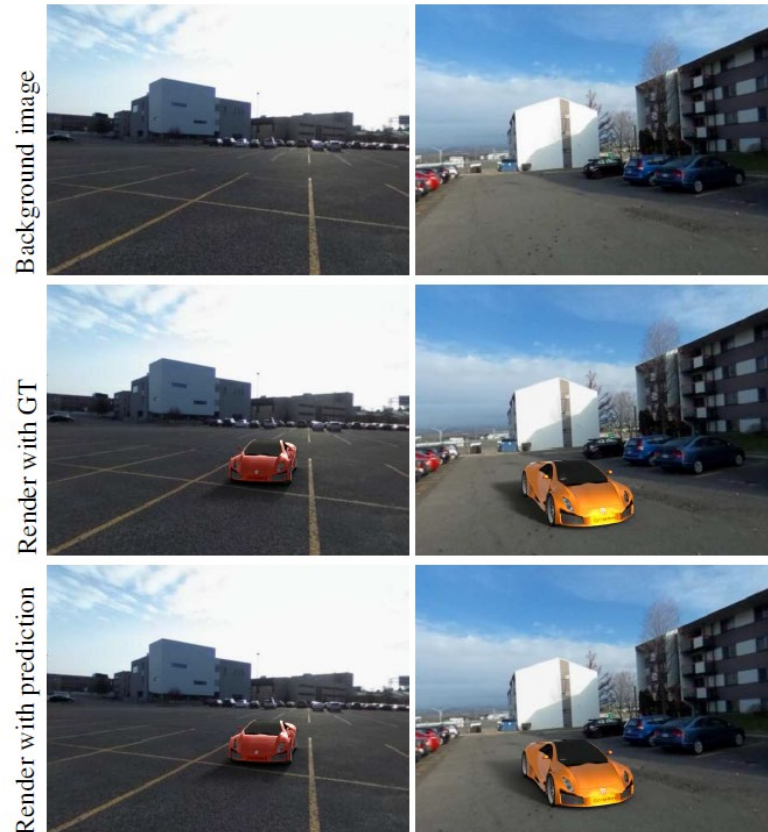
Learning High Dynamic Range from Outdoor Panoramas

Jinsong Zhang Jean-François Lalonde
Université Laval, Québec, Canada

jinsong.zhang.1@ulaval.ca, jflalonde@gel.ulaval.ca

<http://www.jflalonde.ca/projects/learningHDR>

- Regress HDR from one LDR image
- Train on synthetic data
- Limited to outdoor scenes, rotated so that sun is on top



Smarter user assistance

- Handwriting beautification (Zitnick SG'13)
- 3D object modeling (Chen et al. SGA'13)
- 3D object modeling (Kholgade et al. SG'14)

Video and motion

- Video = sequence of images
 - Track points → optical flow, tracked objects, 3D reconstruction
 - Find coherent space-time regions → segmentation
 - Recognizing actions and events
- Examples:
 - Point tracking for structure-from-motion
 - Boujou 1
 - Facial transfer: Xu et al. SG2014

Scene understanding

Interpret image in terms of scene categories, objects, surfaces, interactions, goals, etc.



- Remove the guy lying down (Alyosha)
- Make the woman dance or the guy get up
- Fill in the window with bricks
- Find me images with only Alyosha and Pietro

Scene understanding

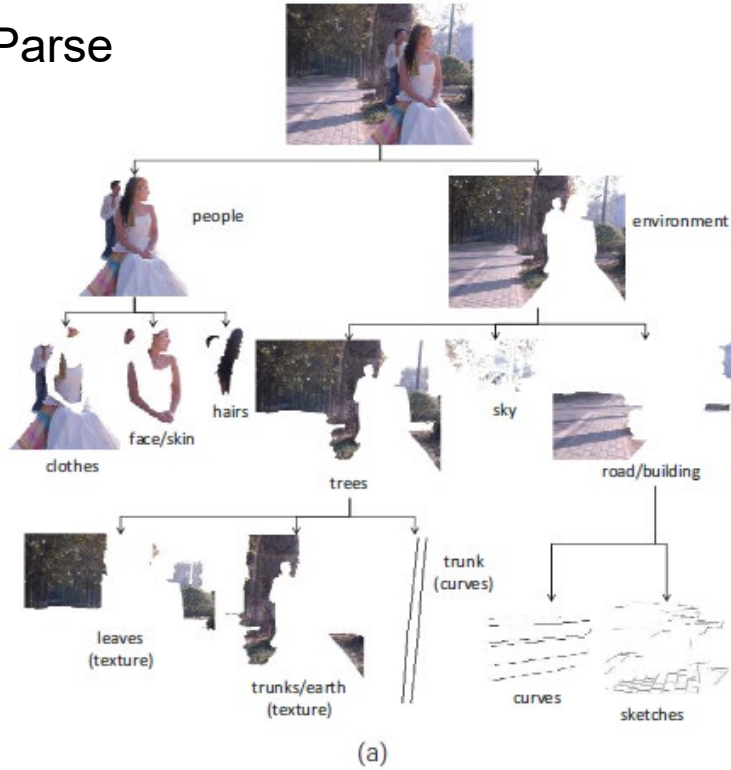
- Mostly unsolved, but we're getting there (especially for graphics purposes)
- Examples
 - “From Image Parsing to Painterly Rendering” (Zeng et al. 2010)
 - “Sketch2Photo: Internet Image Montage” (Chen et al. 2009)
 - Editing via scene attributes (Laffont et al. 2014)

Image Parsing to Painterly Rendering



Image Parsing to Painterly Rendering

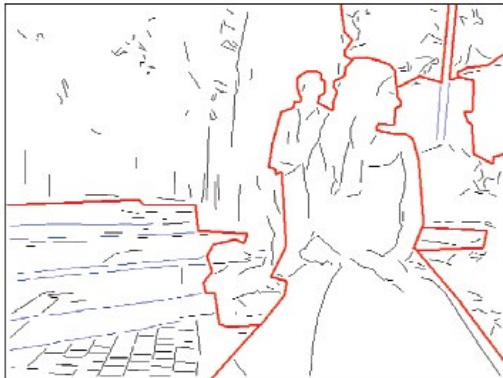
Parse



Brush Strokes



Sketch



Brush Orientations



Image Parsing to Painterly Rendering



Image Parsing to Painterly Rendering



More examples

- Sketch2photo: <http://www.youtube.com/watch?v=dW1Epl2LdFM>
- Animating still photographs



**Animating Pictures
with Stochastic
Motion Textures**

Modeling humans

- Estimating pose and shape
 - <http://clothingparsing.com/>
 - Parselets (Dong et al., ICCV 2013)



- Motion capture
- 3D face from image (Kemelmacher ICCV'13)

Better and simpler 3D reconstruction

MobileFusion (2015): [https://youtu.be/8M -ISYqACo](https://youtu.be/8M-ISYqACo)