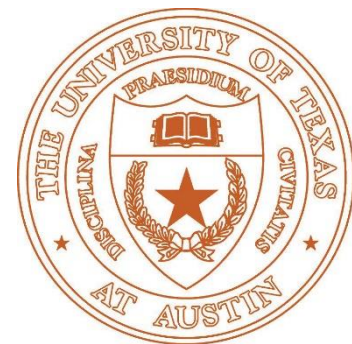


Image Based Reconstruction I

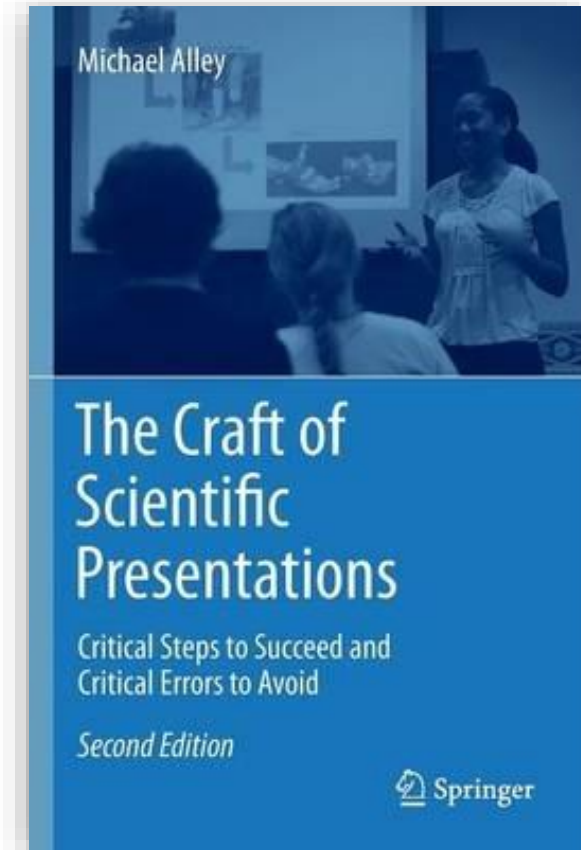
Qixing Huang
Jan. 31th 2017



A Couple of Words on Paper Presentations

- Four components:
 - Motivation
 - Technical Merit
 - Results
 - Broader Impact
- Paper Strength/Weakness
- Read relevant papers as well

Making Presentations



Tools We Will Utilize

- Robust Norms
- MRF Inference
- Continuous Optimization
 - Newton method

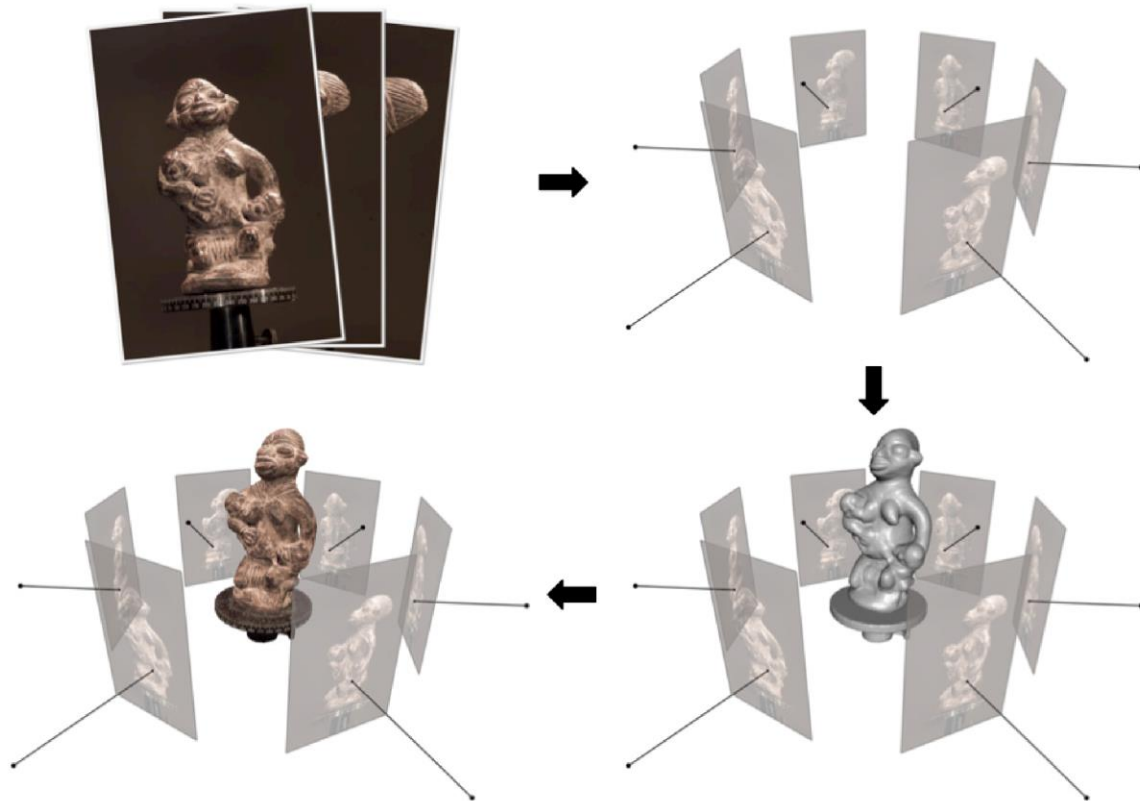
Geometry Reconstruction from Images

3D Reconstruction of a typical
medium-size city

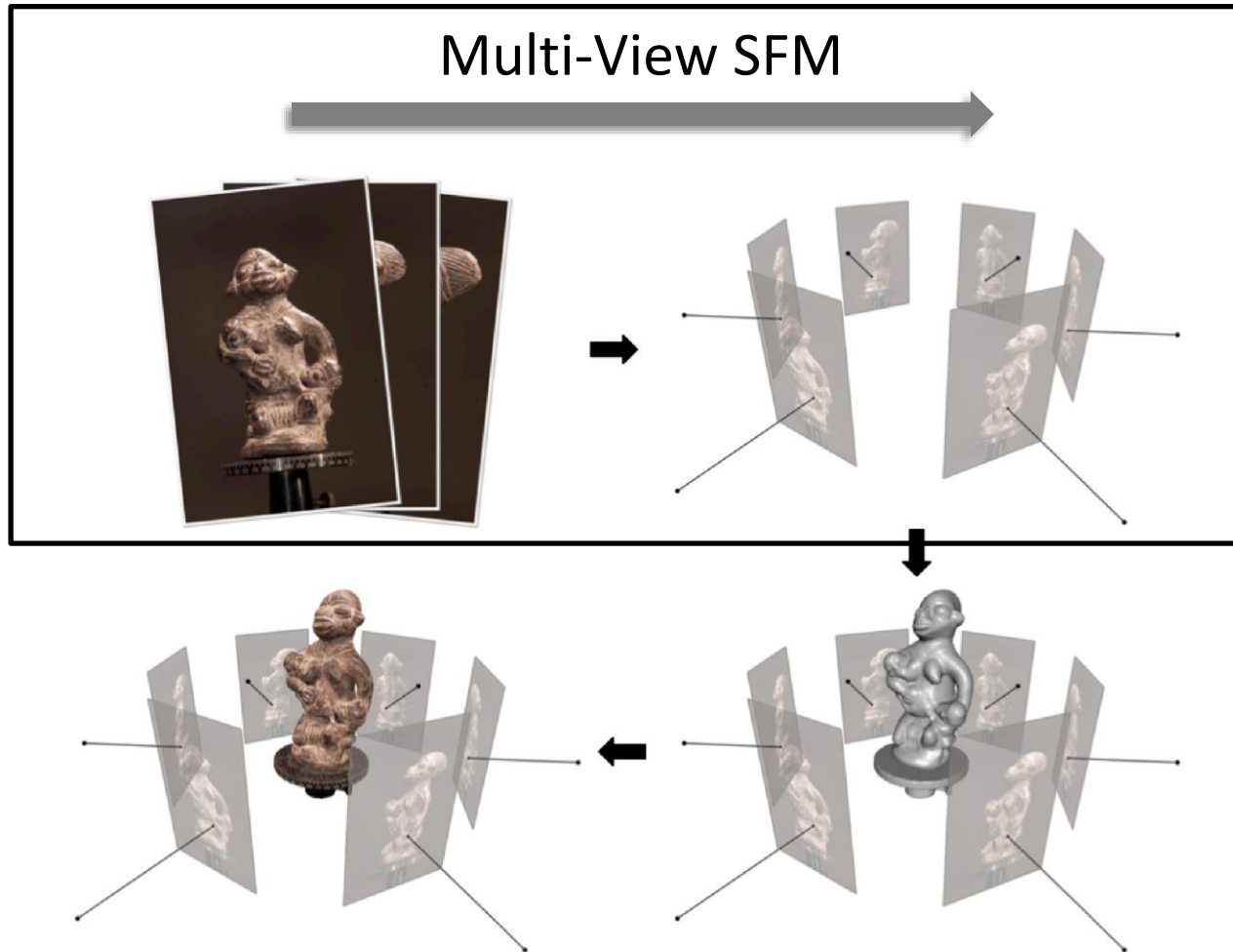
~60,000 images of 50 megapixels

Reconstructed fully automatically in
7 days by 12 servers

Image-Based Geometry Reconstruction Pipeline



This Lecture: Multi-View SFM



SFM Outputs Cameras + (Sparse) Point clouds

[Crandall et al. 13]

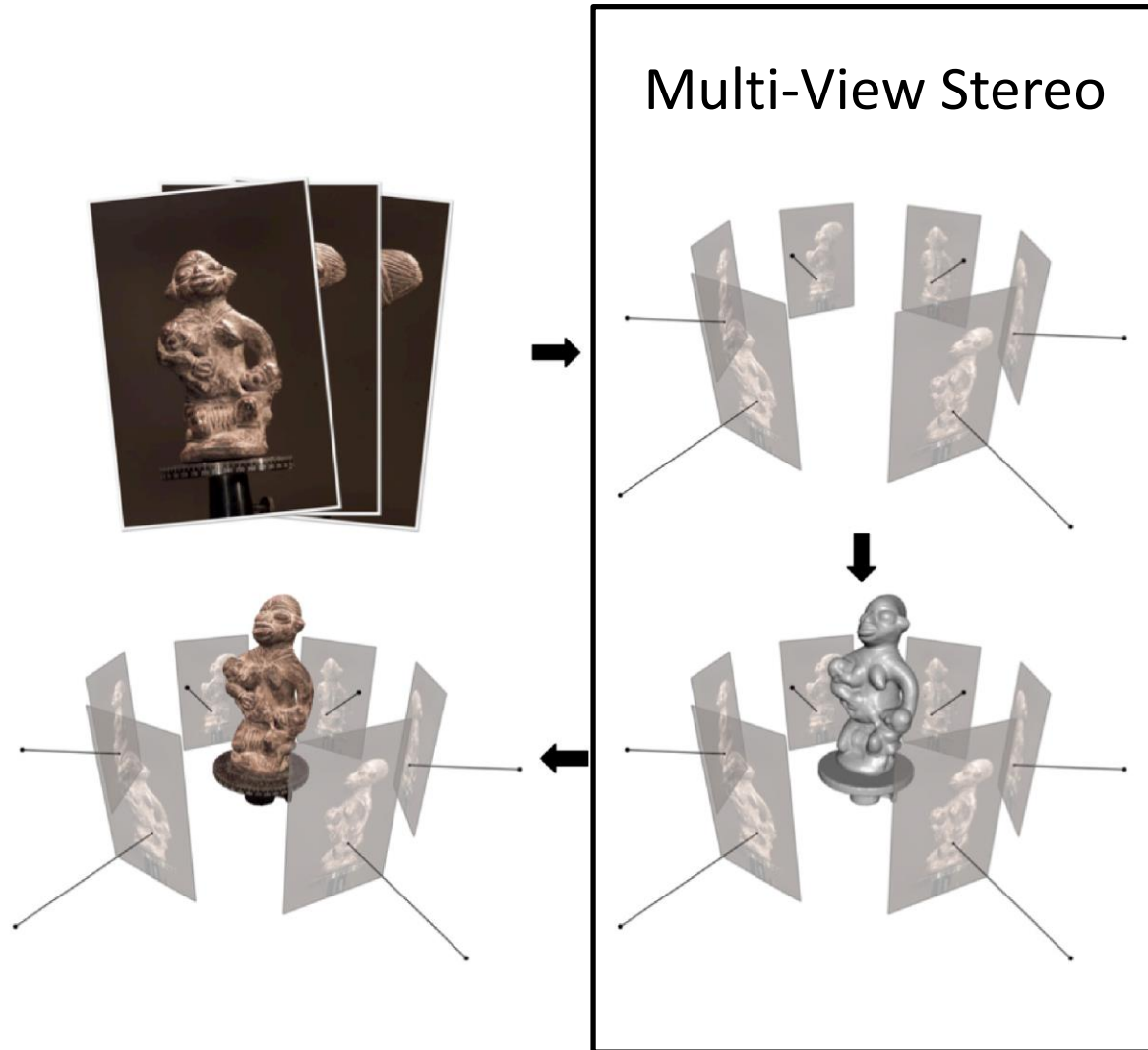


Quad

Reconstructed images: 5,233

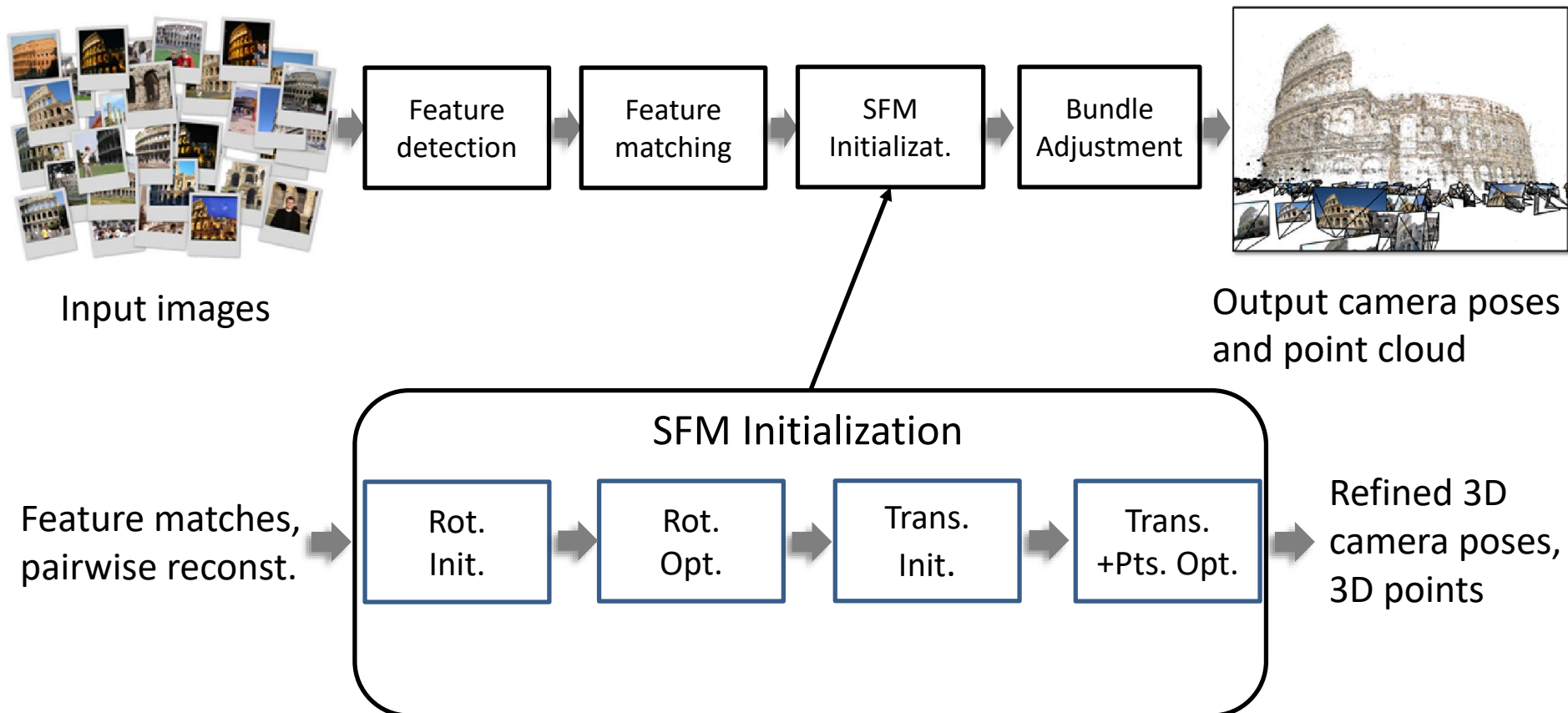
Edges in MRF: 995,734

Next Lecture: Multi-View Stereo

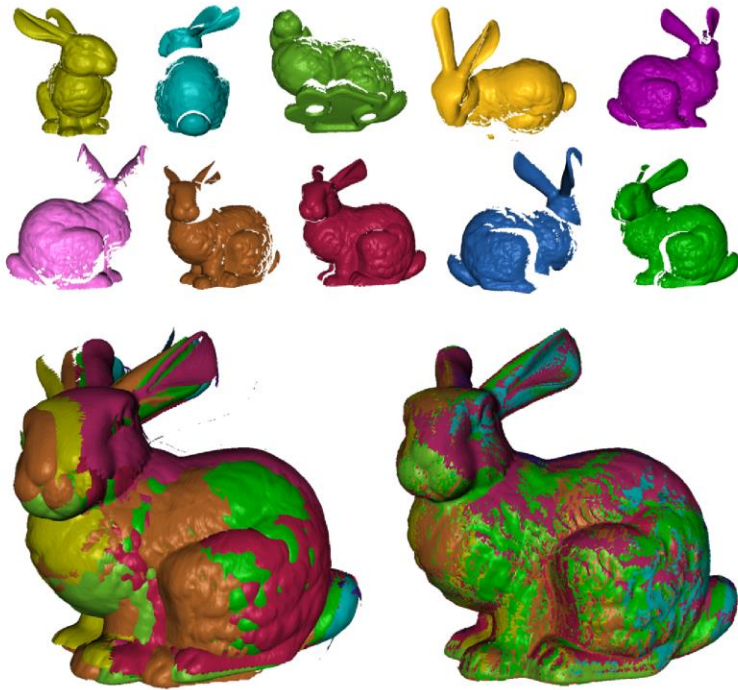


Multi-View SFM Pipeline

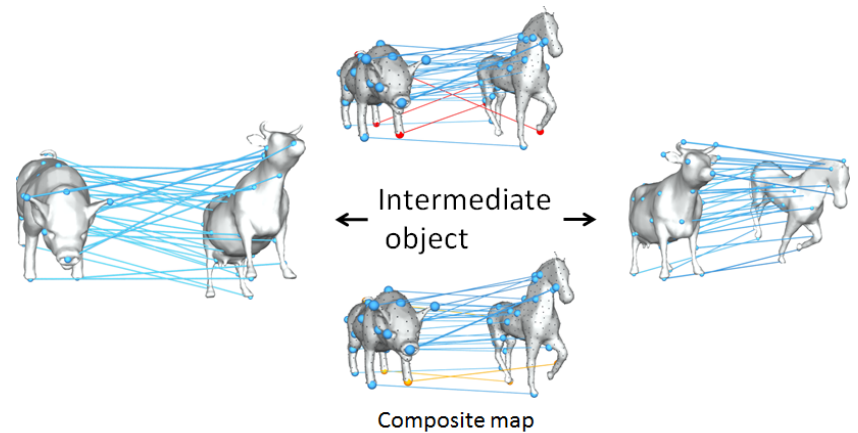
[Crandall et al. 13]



Similar Problems



Scan Alignment
[Gelfand et al. 05]



Data-Driven Map Computation
[Huang et al. 13]

Image Features

SIFT Features [Lowe, IJCV 2004]



Pairwise Image Matching

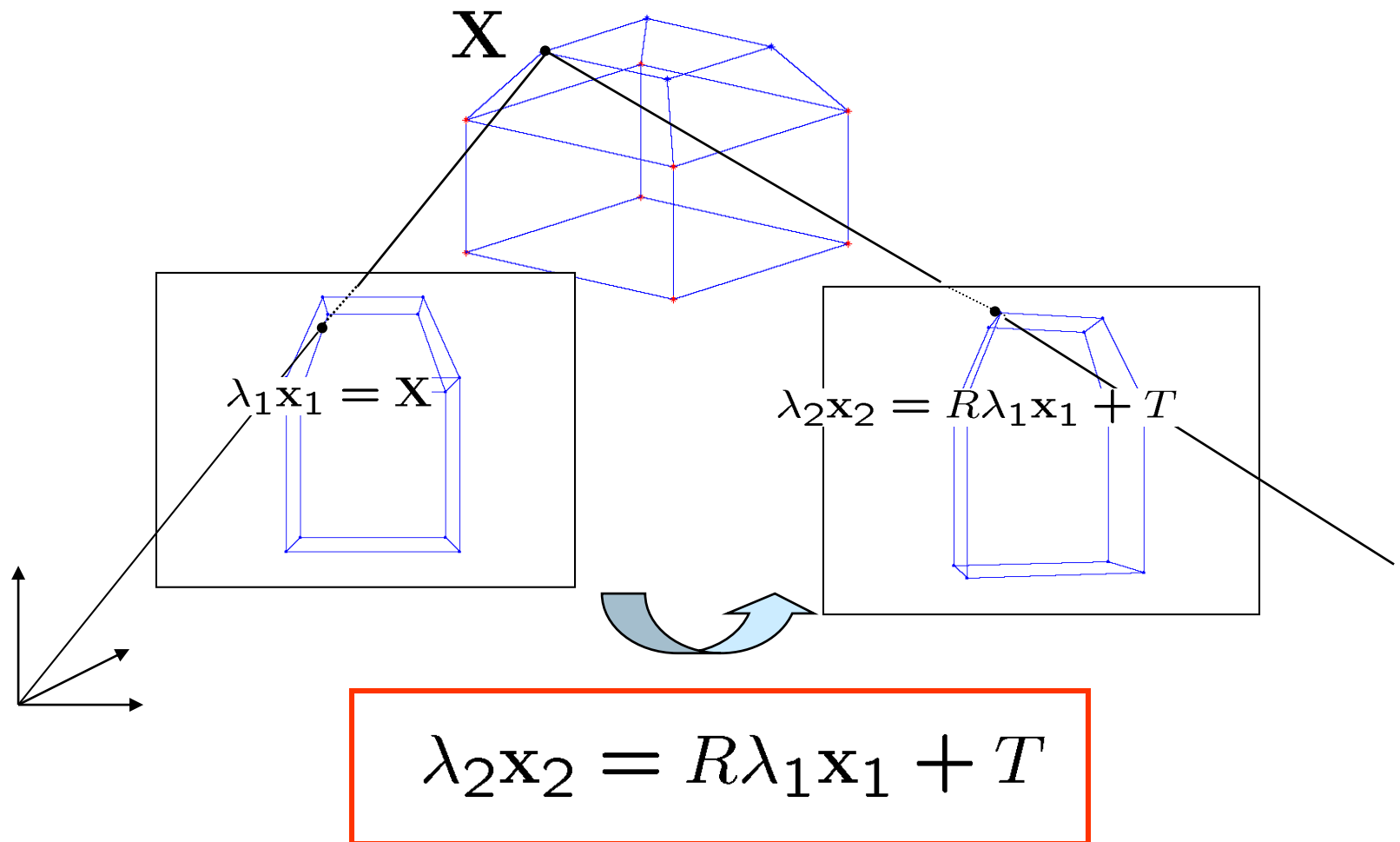
Goal



Given two views of the scene
recover the unknown
relative camera pose

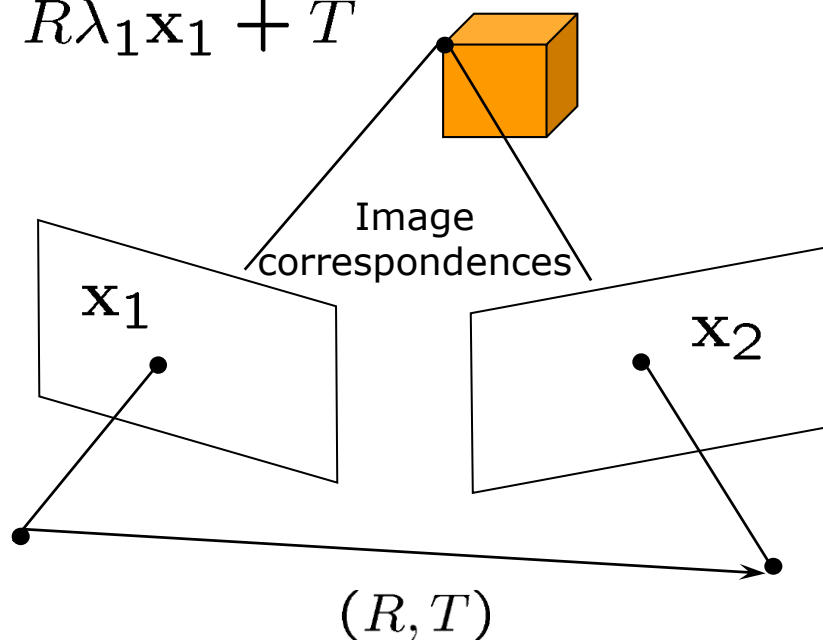
Assume we know the intrinsic camera parameters
Five parameters to optimize

Rigid Body Motion --- Two Views



Epipolar Geometry

$$\lambda_2 \mathbf{x}_2 = R\lambda_1 \mathbf{x}_1 + T$$



- Algebraic Elimination of Depth [Longuet-Higgins '81]:

$$\mathbf{x}_2^T \underbrace{\hat{T}R}_E \mathbf{x}_1 = 0$$

- Essential matrix $E = \hat{T}R$

Nister's Five-Point Method

$$\tilde{q}^\top \tilde{E} = 0$$

$$\tilde{q} \equiv [q_1 q'_1 \quad q_2 q'_1 \quad q_3 q'_1 \quad q_1 q'_2 \quad q_2 q'_2 \quad q_3 q'_2 \quad q_1 q'_3 \quad q_2 q'_3 \quad q_3 q'_3]^\top$$

$$\tilde{E} \equiv [E_{11} \quad E_{12} \quad E_{13} \quad E_{21} \quad E_{22} \quad E_{23} \quad E_{31} \quad E_{32} \quad E_{33}]^\top$$



$$E = xX + yY + zZ + wW$$



$$EE^\top E - \frac{1}{2} \text{trace}(EE^\top) E = 0$$



$$E$$



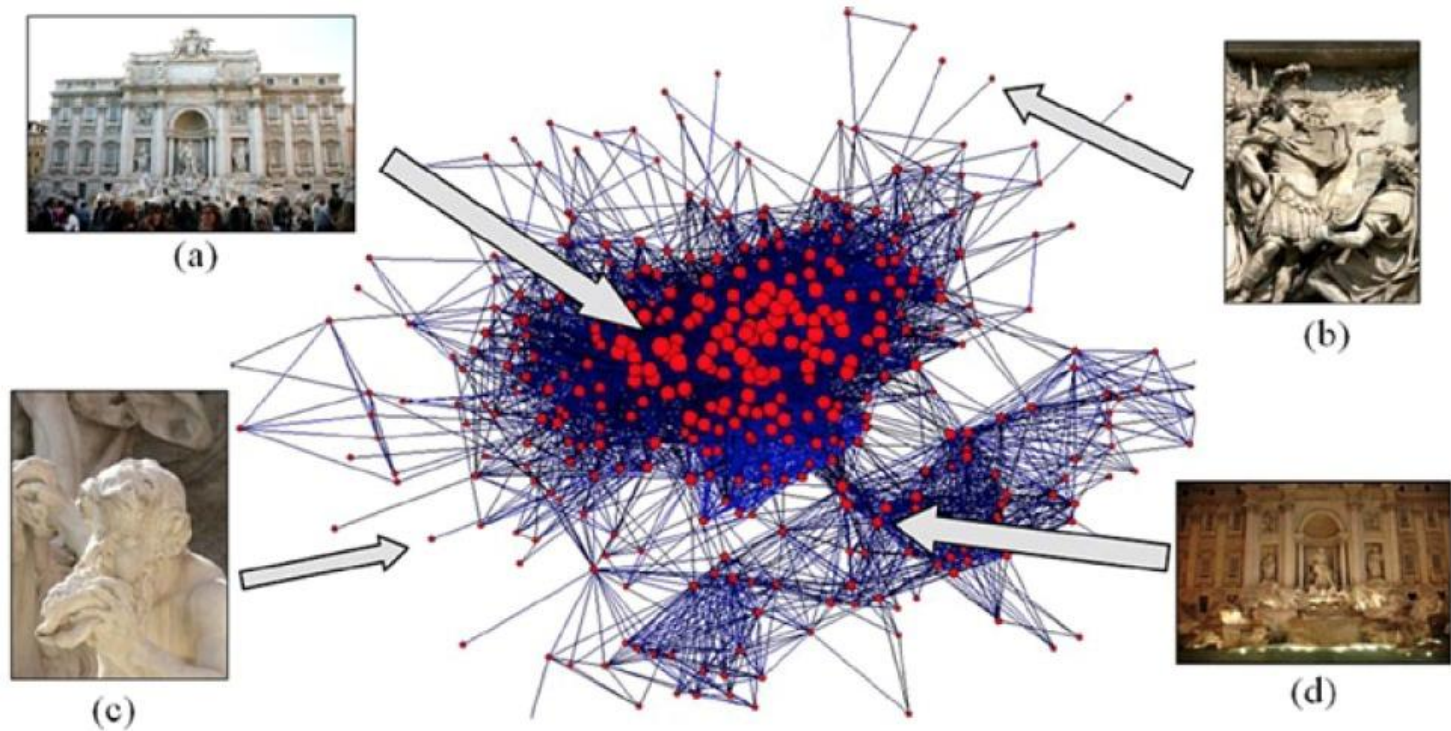
$$R \quad T$$

RANSAC [Fischler and Bolles' 81]

- Pick five feature matches
- Estimate the essential matrix and count the matched SIFT features
- Return the one with the most matched SIFT features

Which Pairs of Images to Match?

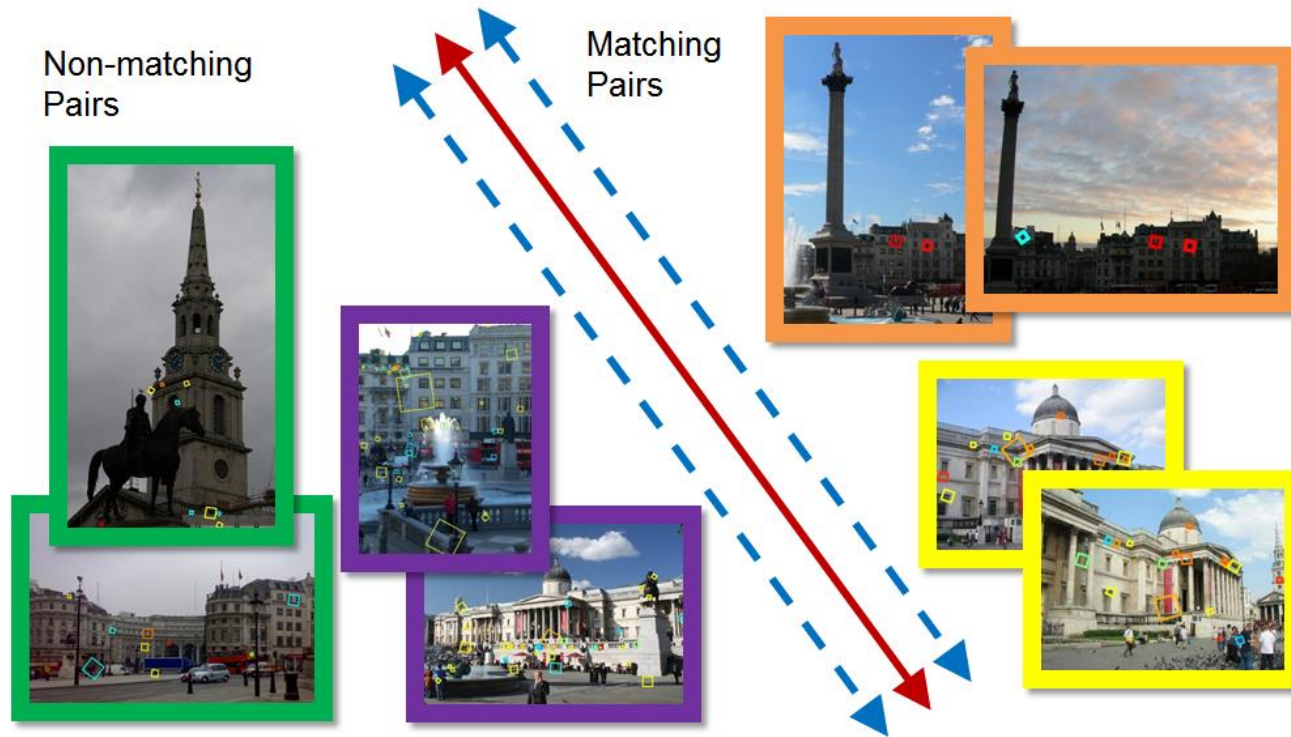
Nearest Neighbors in Image Descriptors (e.g., GIST and HOG)



Only works for images that significantly overlap

Train a Classifier to Differentiate Good/Bad Matches

[Cao and Snavely' 12]



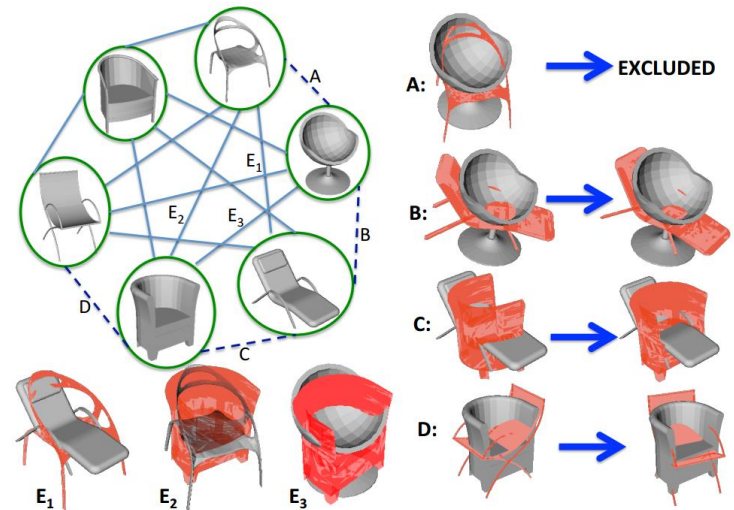
Iterative Matching/Learning

- Start from matching image descriptors
- Verification via image matching
- Train a classifier to find more potential image pairs and iterate

Graph Connectivity Optimization --- maximizing $\lambda_2(G)$



Imageweb [Heath et al 10]



Fuzzy correspondences
on shapes [Kim et al 12]

Multi-View Pose Estimation

What We Have So Far

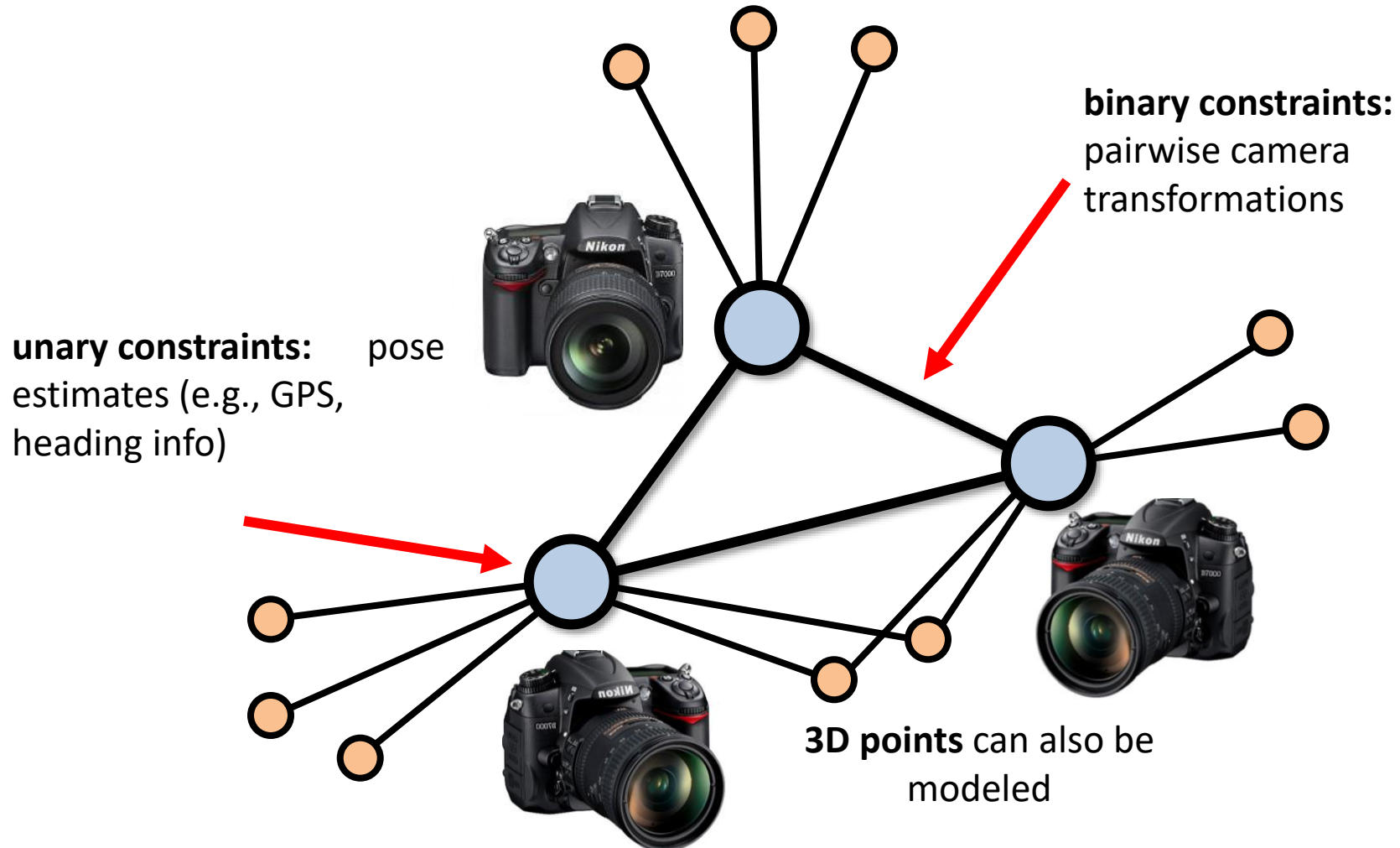
- A graph of images
- Along each edge
 - Noisy relative poses
 - Matched SIFT features

Three Approaches of Multi-View Pose Estimation (Goal: Remove Bad Matches)

- Combinatorial Optimization
- Convex/Nonconvex Optimizations
- MRF-Based Formulation

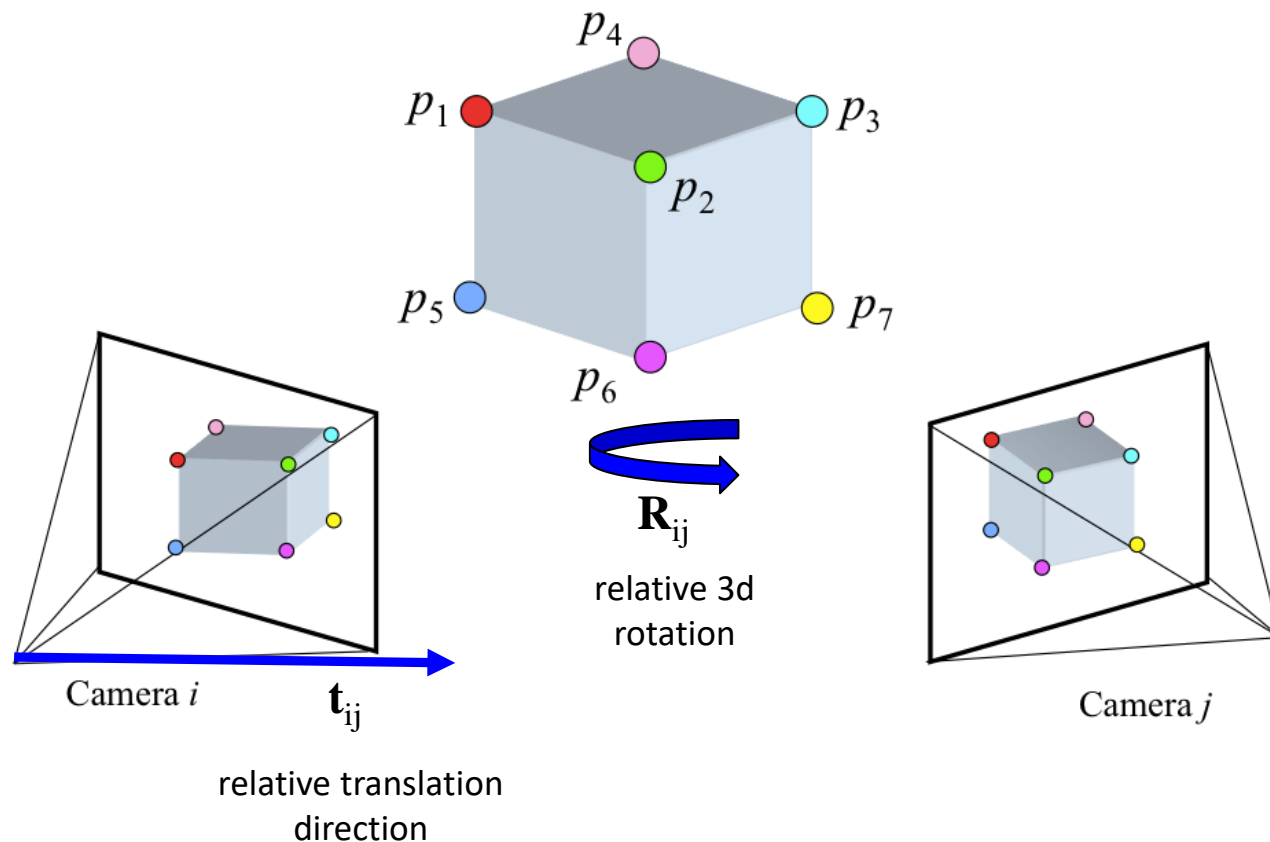
The MRF model

- Input: set of images with correspondence

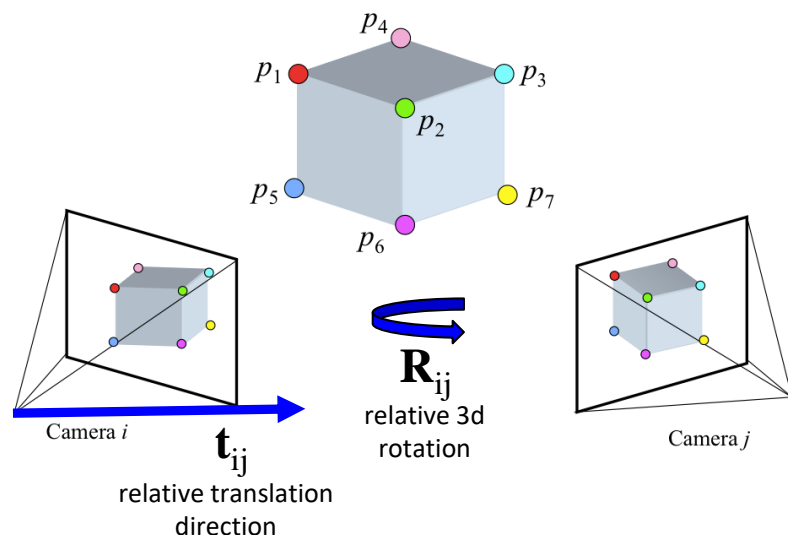


Constraints on camera pairs

- Compute relative pose between camera pairs using 2-frame SfM [Nister04]



Constraints on camera pairs



- Find absolute camera poses $(\mathbf{R}_i, \mathbf{t}_i)$ and $(\mathbf{R}_j, \mathbf{t}_j)$ that agree with these pairwise estimates:

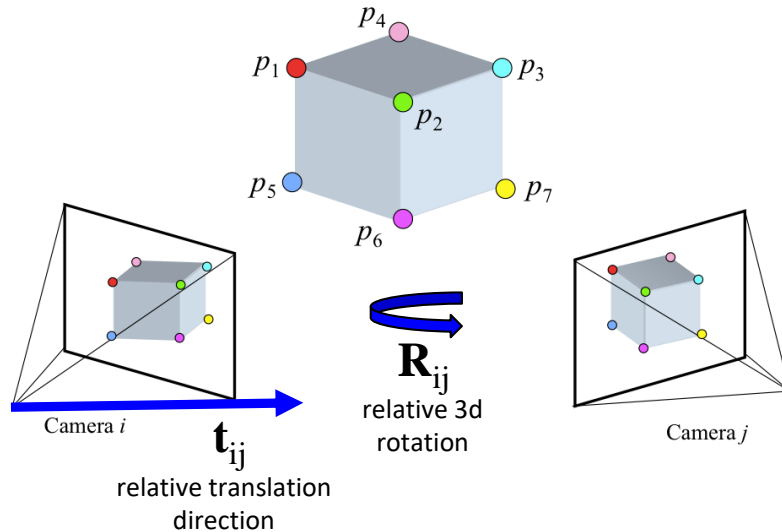
$$\mathbf{R}_{ij} = \mathbf{R}_i^\top \mathbf{R}_j$$

rotation consistency

$$\lambda_{ij} \mathbf{t}_{ij} = \mathbf{R}_i^\top (\mathbf{t}_j - \mathbf{t}_i)$$

translation direction consistency

Constraints on camera pairs



- Define robustified error functions to use as pairwise potentials:

$$d^{\mathbf{R}}(\mathbf{R}_{ij}, \mathbf{R}_i^{\top} \mathbf{R}_j)$$

$$d^{\mathbf{R}}(\mathbf{R}_a, \mathbf{R}_b) = \rho_R(\|\mathbf{R}_a - \mathbf{R}_b\|)$$

rotation consistency

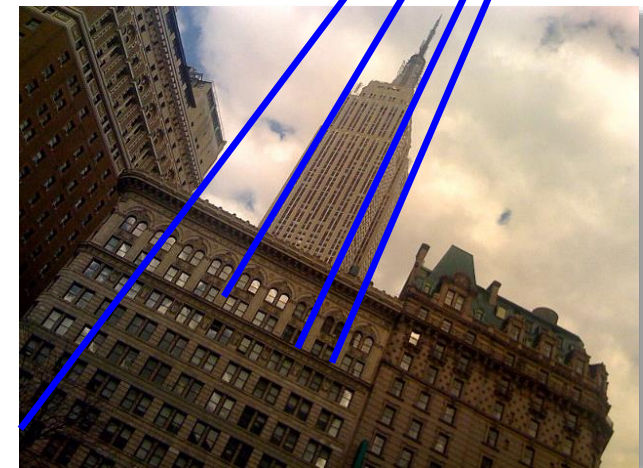
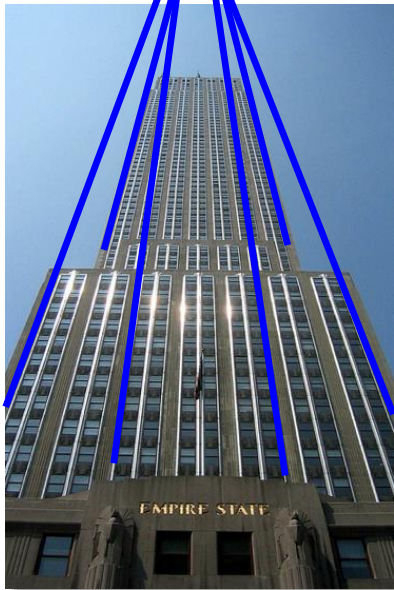
$$d^{\mathbf{T}}(\mathbf{t}_j - \mathbf{t}_i, \mathbf{R}_i \mathbf{t}_{ij})$$

$$d^{\mathbf{T}}(\mathbf{t}_a, \mathbf{t}_b) = \rho(\text{angleof}(\mathbf{t}_a, \mathbf{t}_b))$$

translation direction consistency

Prior pose information

- *Noisy* absolute pose info for some cameras
 - 2D positions from geotags (GPS coordinates)
 - Orientations (tilt & twist angles) from vanishing point detection [Sinha10]



Overall optimization problem

- Given pairwise and unary pose constraints, solve for absolute camera poses simultaneously
 - for n cameras, estimate

$$\mathcal{R} = (\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_n) \quad \text{and} \quad \mathcal{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n)$$

so as to minimize total error over the entire graph

$$D^{\mathbf{R}}(\mathcal{R}) = \sum_{e_{ij} \in E_C} \boxed{d^{\mathbf{R}}(\mathbf{R}_{ij}, \mathbf{R}_i^{\top} \mathbf{R}_j)} + \alpha_1 \sum_{I_i \in \mathcal{I}} \boxed{d_i^{\mathbf{O}}(\mathbf{R}_i)}$$

pairwise rotation consistency unary rotation consistency

$$D^{\mathbf{T}}(\mathcal{T}, \mathcal{R}) = \sum_{e_{ij} \in E_C} \boxed{d^{\mathbf{T}}(\mathbf{t}_j - \mathbf{t}_i, \mathbf{R}_i \mathbf{t}_{ij})} + \alpha_2 \sum_{I_i \in \mathcal{I}} \boxed{d_i^{\mathbf{G}}(\mathbf{t}_i)}$$

pairwise translation consistency unary translation consistency

MRF Inference

- Convert continuous optimization into a labeling problem:

$$E(\mathbf{x}) = \sum_i f_i(x_i) + \sum_{(i,j) \in \mathcal{E}} f_{ij}(x_i, x_j)$$

- A well studied problem with efficient solvers

Incorporate Points

- Point tracks --- interest points across multiple images that have similar SIFT descriptors
- Select point tracks that cover each camera-camera edge five times and each image ten times
- Relation between 3D location and image coordinates:

$$\mu_{ik} \mathbf{x}_{ik} = \mathbf{K}_i \mathbf{R}_i (\mathbf{X}_k - \mathbf{t}_i)$$

Intrinsic camera parameters

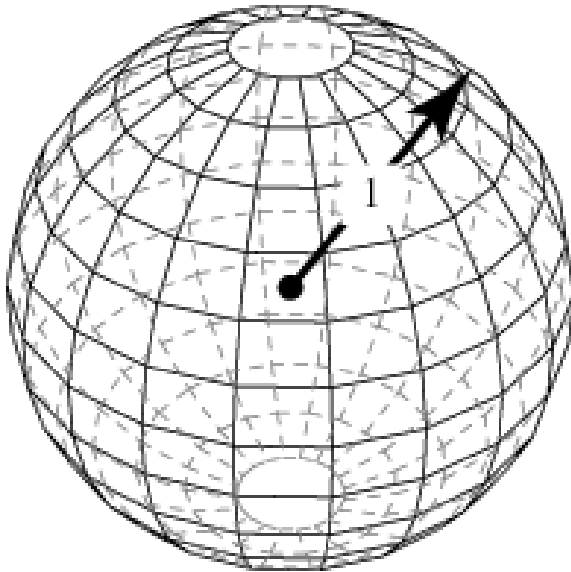
Solving the MRF

- Use discrete loopy belief propagation [Pearl88]
 - Up to 1,000,000 nodes (cameras and points)
 - Up to 5,000,000 edges (constraints between cameras and points)
 - 6-dimensional label space for cameras (3-dimensional for points)

Solving the MRF

- Reduce 6-dimensional label space by...
 - Solving for rotations & translations independently [Martinec07], [Sim06], [Sinha08]
 - Assuming camera twist angles are near 0
 - Initially solving for 2D camera positions
- Speed up BP by...
 - Using a parallel implementation on a cluster
 - Using distance transforms (aka min convolutions) to compute BP messages in $O(L)$ time in # of labels (instead of $O(L^2)$) [Felzenszwalb04]

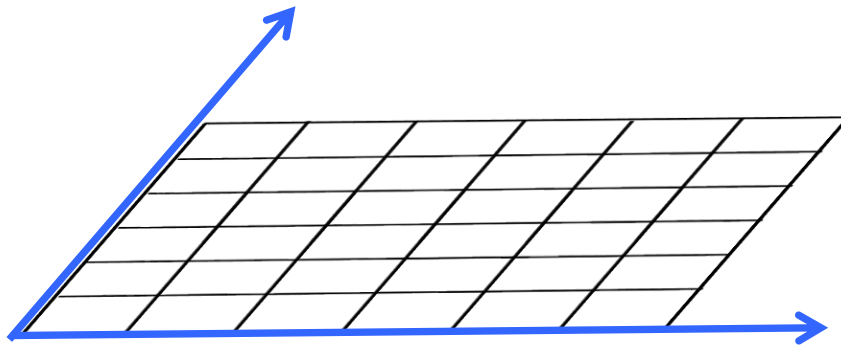
Discrete BP: Rotations



- Parameterize viewing directions as points on unit sphere
 - Discretize into $10 \times 10 \times 10 = 1,000$ possible labels
 - Measure rotational errors as robust Euclidean distances on sphere (to allow use of distance transform)


Discrete BP: Translations

- Parameterize positions as 2D points in plane
 - Use approximation to error function
(to allow use of distance transforms)
 - Discretize into up to $300 \times 300 = 90,000$ labels



Bundle Adjustment

Rotation Optimization

$$D^{\mathbf{R}}(\mathcal{R}) = \sum_{e_{ij} \in E_C} d^{\mathbf{R}}(\mathbf{R}_{ij}, \mathbf{R}_i^{\top} \mathbf{R}_j)$$
$$d^{\mathbf{R}}(\mathbf{R}_a, \mathbf{R}_b) = \rho_R(\|\mathbf{R}_a - \mathbf{R}_b\|)$$


The diagram shows a blue arrow pointing from the x^2 term to the ρ_R function in the second equation, indicating that ρ_R is a quadratic loss function.

Using quadratic loss after removing outlier rotations

Fix one image (or use geotags which provide pose priors)

Nonlinear Least Squares



Gauss-Newton Method

Gauss-Newton Method

- The Gauss–Newton algorithm is a method used to solve non-linear least squares problems

$$f(x) \equiv \frac{1}{2} \Delta z(x)^T W \Delta z(x)$$

$$g \equiv \frac{df}{dx} = \Delta z^T W J$$

$$H \equiv \frac{d^2 f}{dx^2} = J^T W J + \sum_i (\Delta z^T W)_i \frac{d^2 z_i}{dx^2}$$

$$H \approx J^T W J, \quad \frac{d^2 z_i}{dx^2} \approx 0$$

$$(J^T W J) \delta x = -J^T W \Delta z$$

Linear Convergence for Practical Problems

Levenberg – Marquardt Heuristic

- The LMA interpolates between the Gauss–Newton algorithm (GNA) and the method of gradient descent
- When far from the minimum it acts as a steepest descent and it performs gauss newton iteration when near to the solution

$$(H + \lambda W) \delta x = -g$$

Translation and Point Optimization

Camera-Camera
Relation:

$$\lambda_{ij} \mathbf{t}_{ij} = \mathbf{R}_i^\top (\mathbf{t}_j - \mathbf{t}_i) \quad \hat{\mathbf{t}}_{ij} = \mathbf{R}_i \mathbf{t}_{ij}$$


Camera-Point
Relation:

$$\mu_{ik} \mathbf{x}_{ik} = \mathbf{K}_i \mathbf{R}_i (\mathbf{X}_k - \mathbf{t}_i) \quad \hat{\mathbf{x}}_{ik} = \mathbf{R}_i^\top \mathbf{K}_i^{-1} \mathbf{x}_{ik}$$

$$D^{\mathbf{T}}(\mathcal{T}) = \alpha_2 \sum_{e_{ij} \in E_C} d^{\mathbf{T}}(\mathbf{t}_j - \mathbf{t}_i, \hat{\mathbf{t}}_{ij}) + d^{\mathbf{T}}(\mathbf{t}_i - \mathbf{t}_j, \hat{\mathbf{t}}_{ji}) +$$


$$\alpha_3 \sum_{e_{ik} \in E_F} d^{\mathbf{T}}(\mathbf{X}_k - \mathbf{t}_i, \hat{\mathbf{x}}_{ik})$$

$$d^{\mathbf{T}}(\mathbf{v}_a, \mathbf{v}_b) = \rho(\text{angleof}(\mathbf{v}_a, \mathbf{v}_b))$$


 x^2


Positional variables are decoupled!

Schur Trick

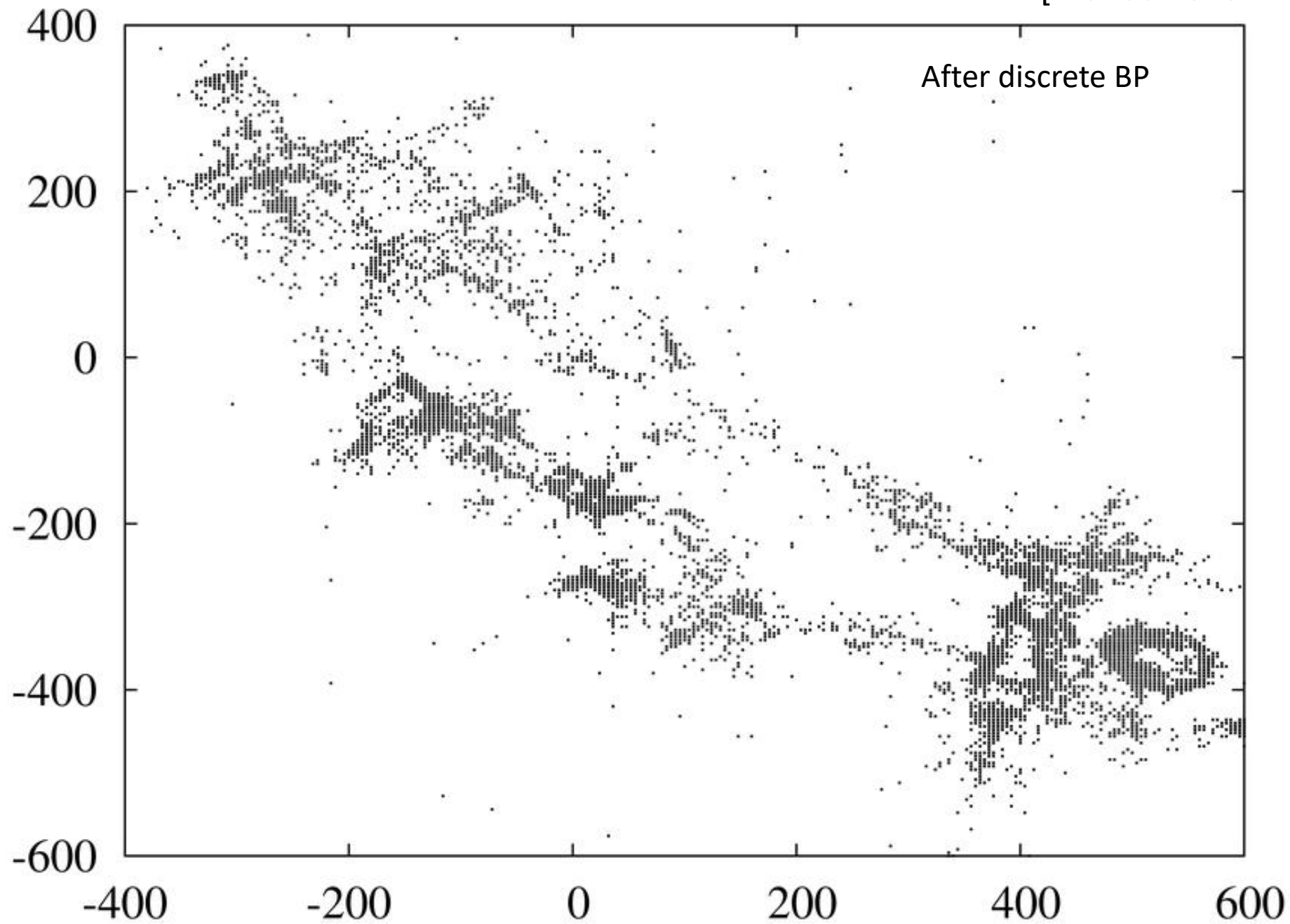
$$\begin{array}{c} \text{Trans.} \\ \text{Poss.} \end{array} \begin{array}{cc} \text{Trans.} & \text{Poss.} \\ \left[\begin{array}{cc} B & E \\ E^\top & C \end{array} \right] & \begin{bmatrix} \Delta y \\ \Delta z \end{bmatrix} = \begin{bmatrix} v \\ w \end{bmatrix} \end{array}$$


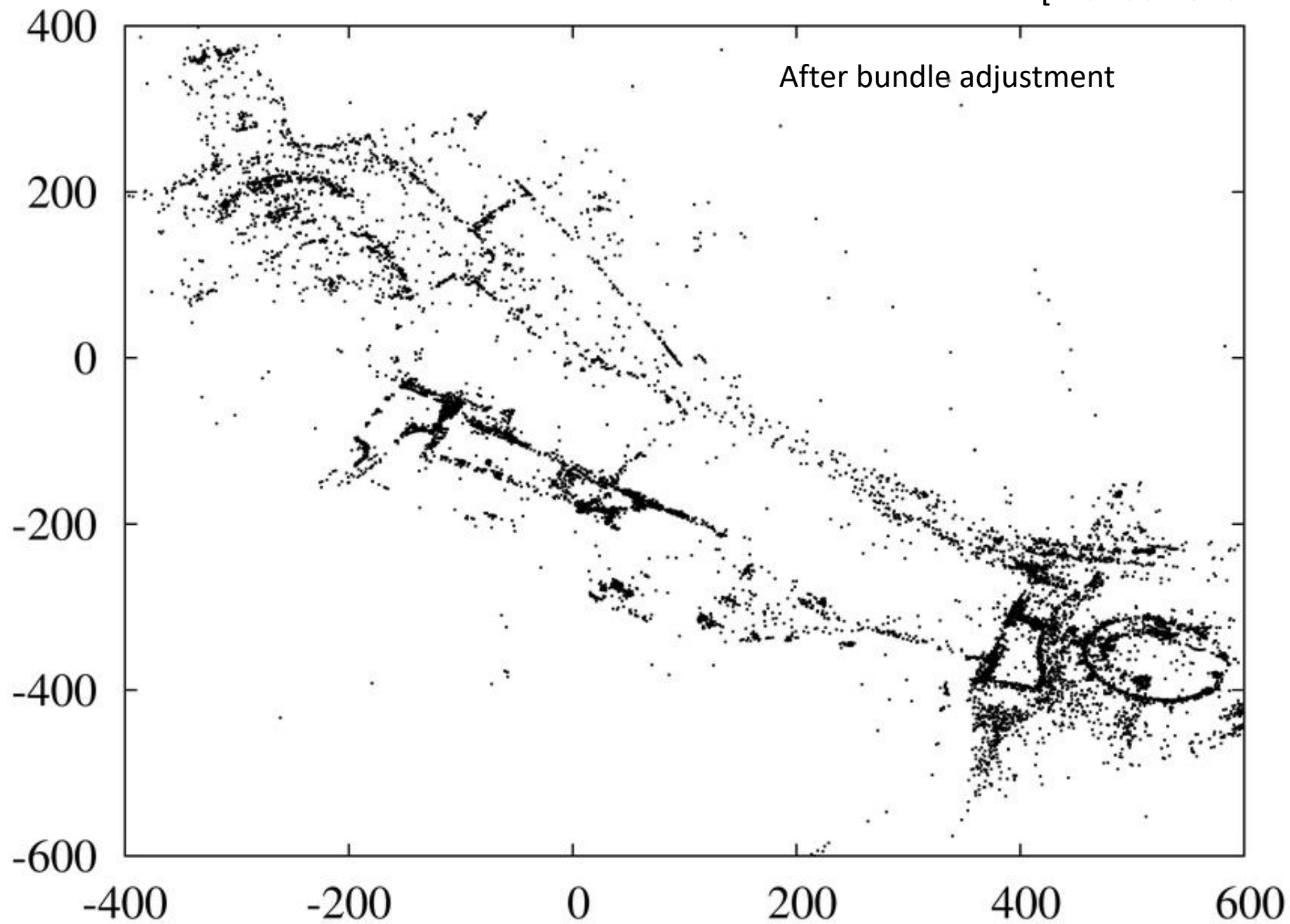
Block Diagonal Matrix

$$\Delta z = C^{-1} (w - E^\top \Delta y)$$

$$\left[B - EC^{-1}E^\top \right] \Delta y = v - EC^{-1}w$$


Small-scale linear system





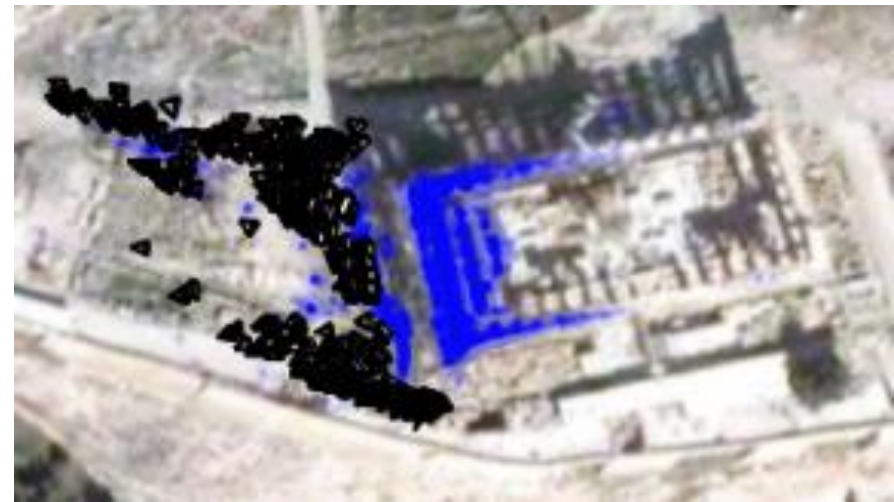
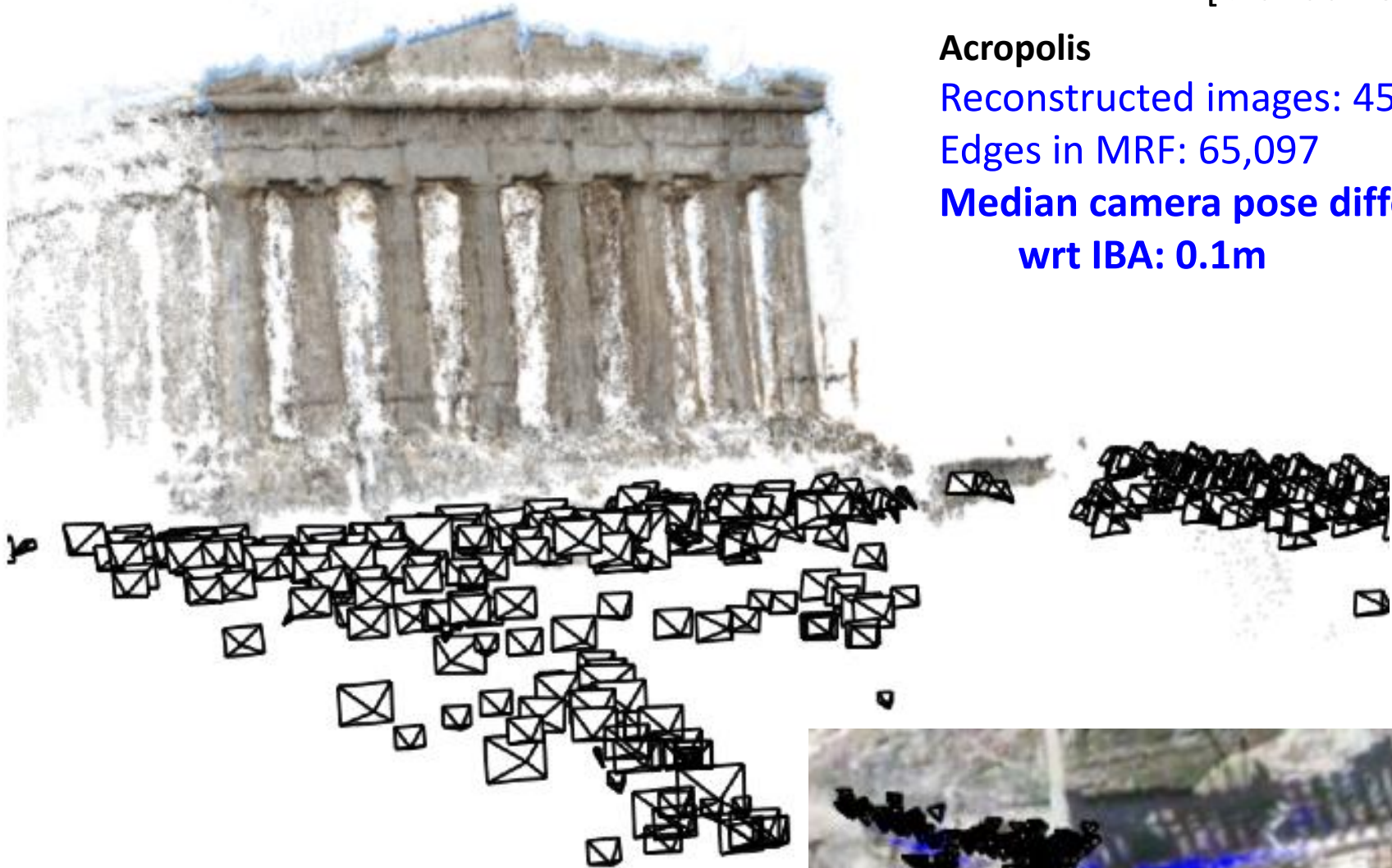
Experimental Results

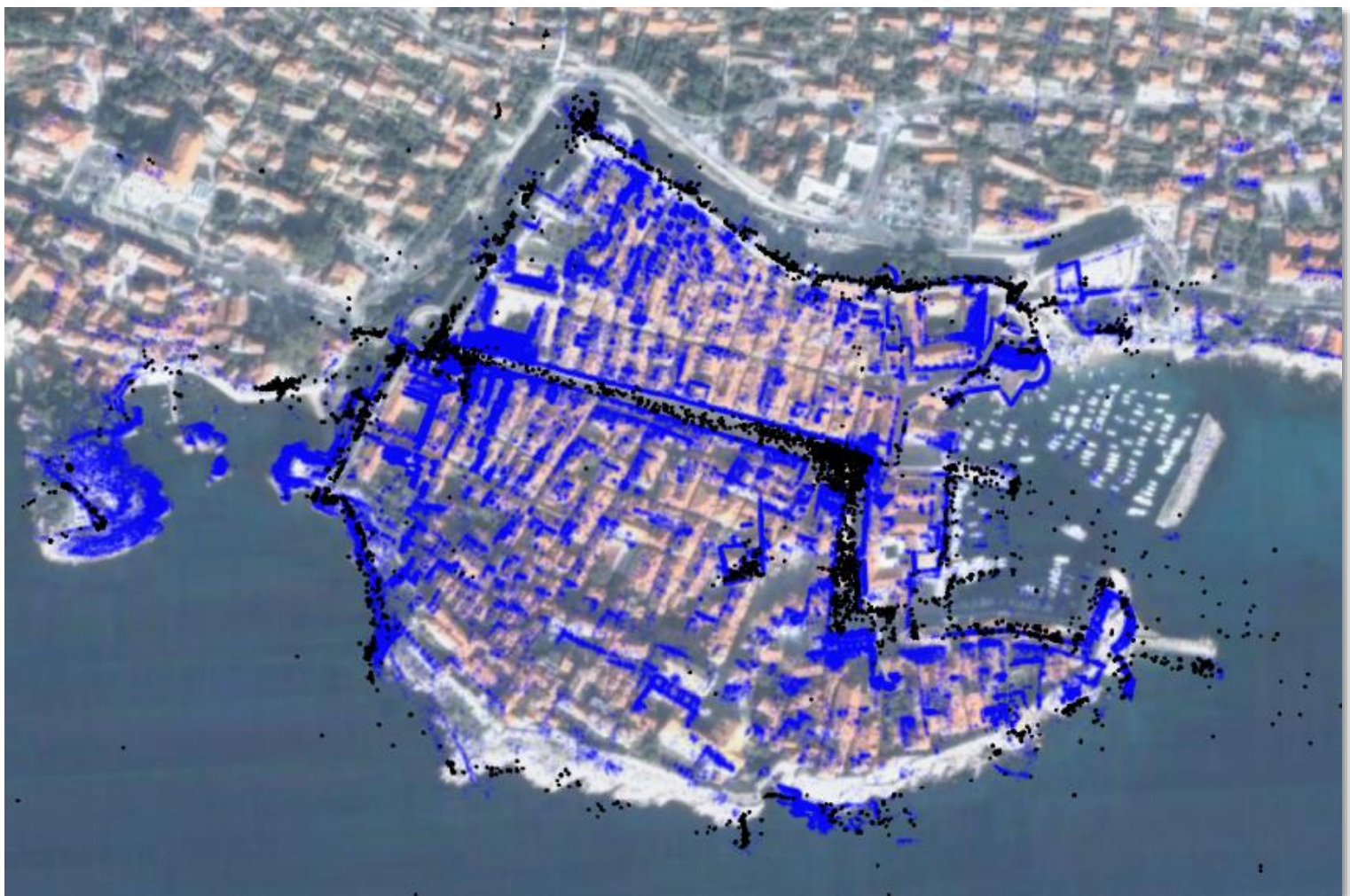
Acropolis

Reconstructed images: 454

Edges in MRF: 65,097

**Median camera pose difference
wrt IBA: 0.1m**





Dubrovnik (Croatia)

Reconstructed images: 6,532

Edges in MRF: 1,835,488

Median camera pose difference wrt IBA: 1.0m

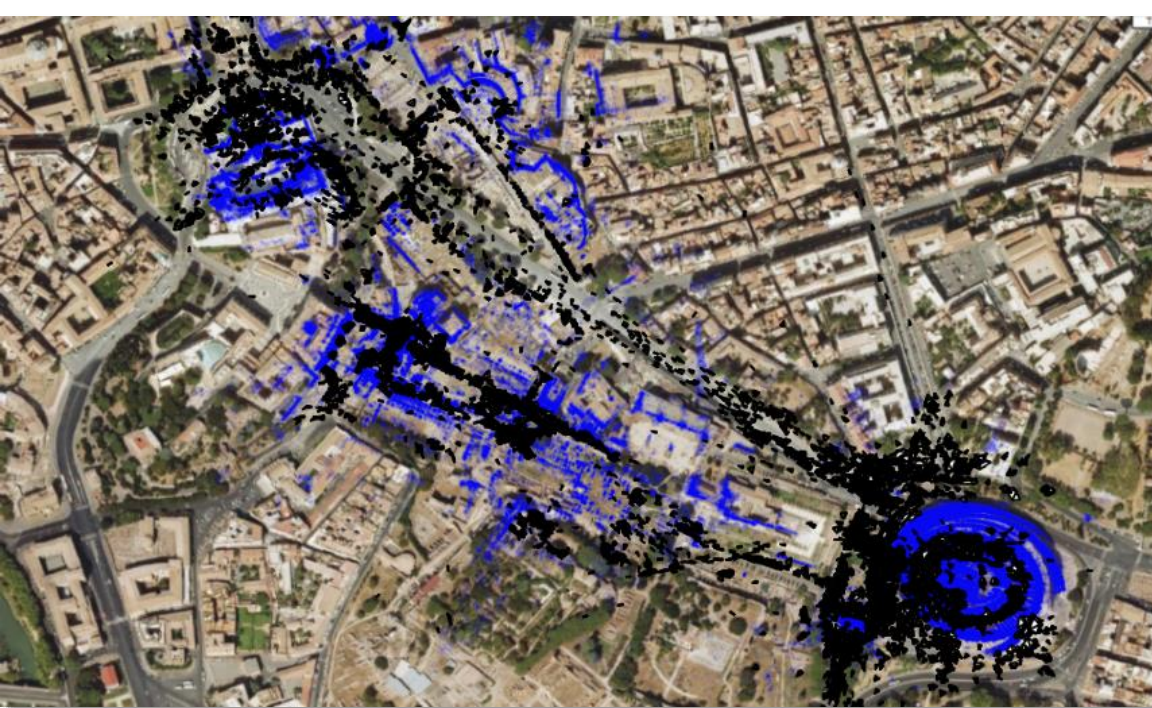
Central Rome

[Crandall et al. 13]

Reconstructed images: 14,754

Edges in MRF: 2,258,416





Central Rome

Reconstructed images: 14,754

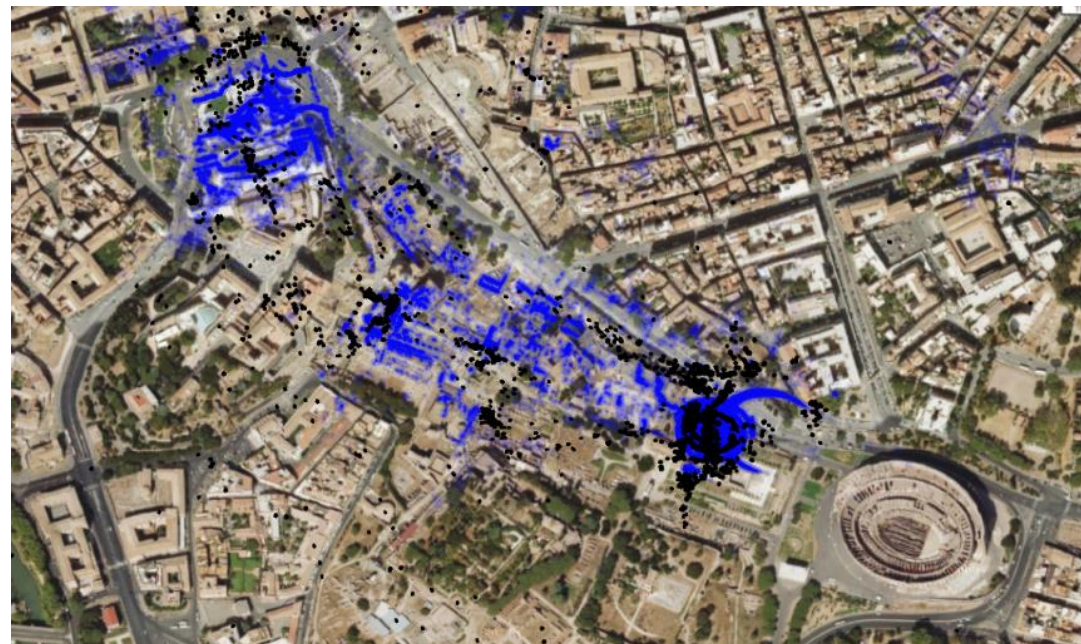
Edges in MRF: 2,258,416

**Median camera pose difference
wrt IBA: 25.0m**

Our result

Incremental
Bundle Adjustment
[Agarwal09]

[Crandall et al. 13]



How can Deep Learning Help?

Pipeline Steps

- Pairwise matching?
- Graph reconstruction?
- Multi-image matching?