# Efficient Global 2D-3D Matching for Camera Localization in a Large-Scale 3D Map

**Liu Liu, Hongdong Li, and Yuchao Dai**

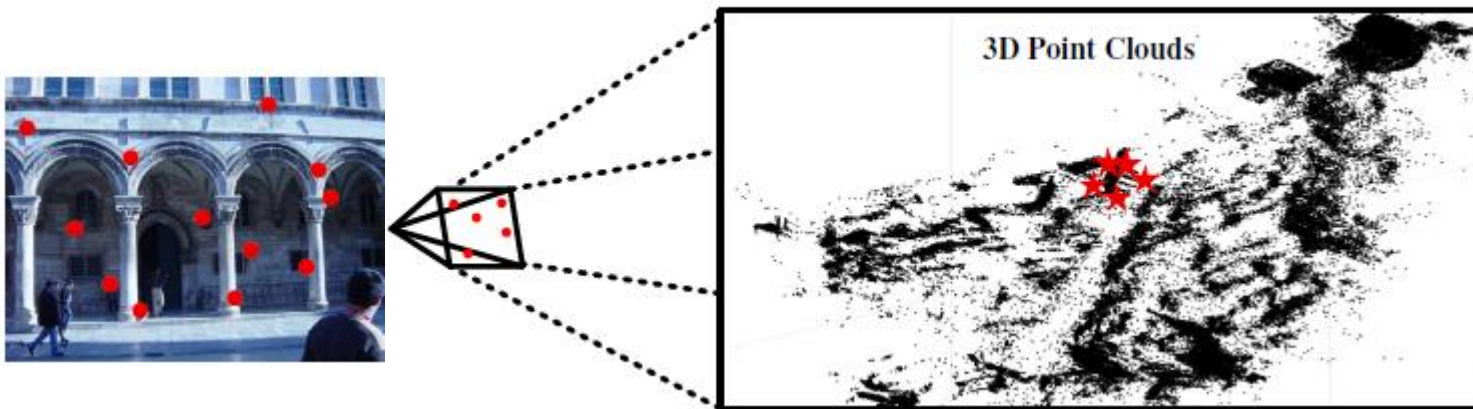*In Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017*

**Speaker: B. Y. Huang**

# Outline

- Research field
- Motivation
- Proposed scheme
- Experiment and comparison
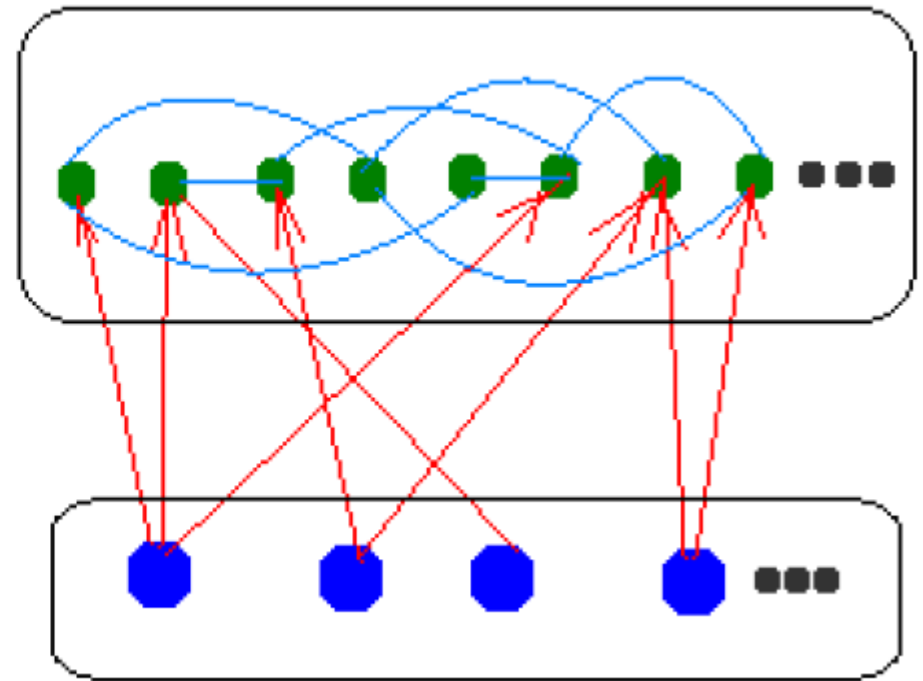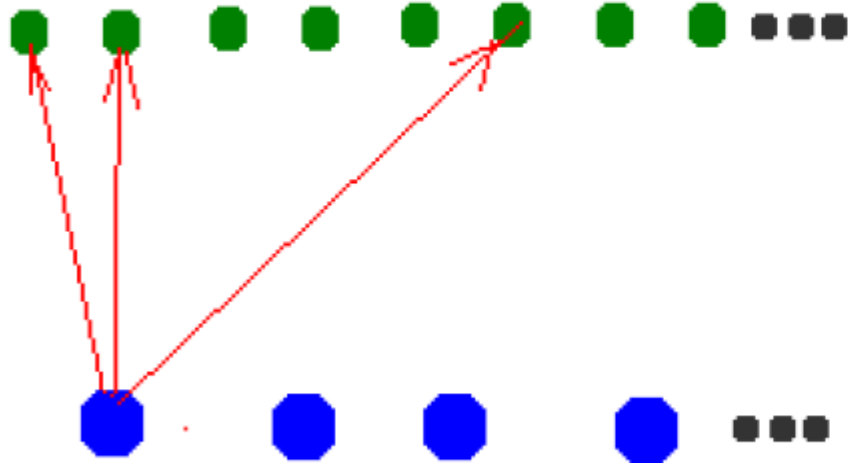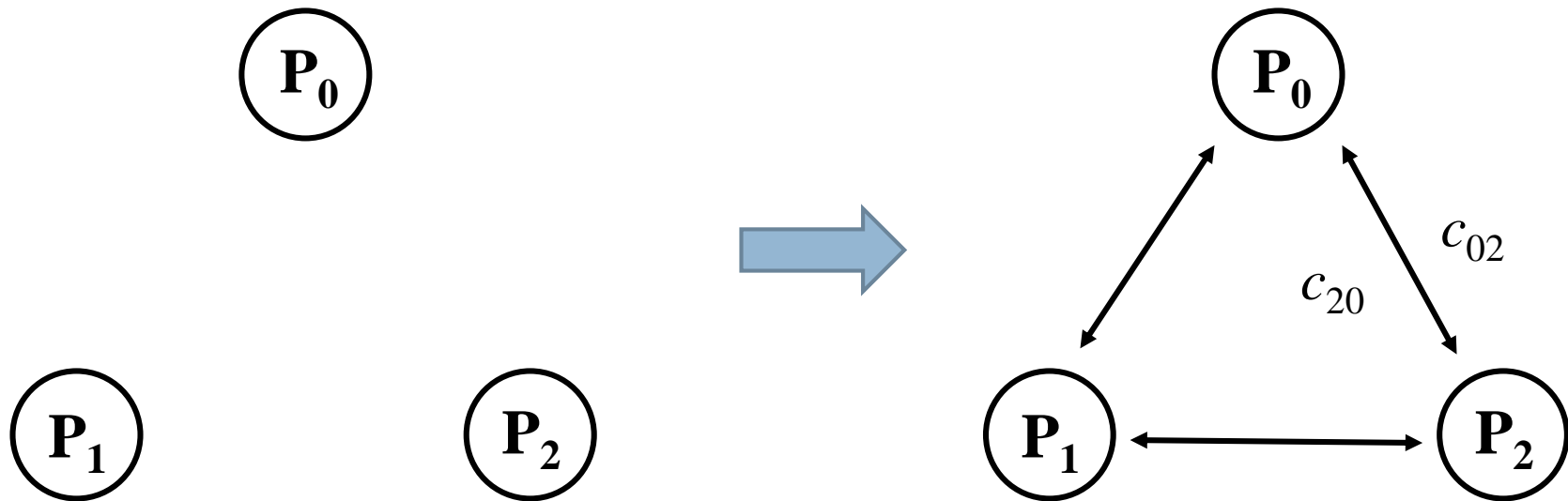- Conclusion

# What?

- Image-based localization

# Why?

- **Large-scale problem → ambiguity**
  - 3D points can be <span style="color:red">visually similar</span> or even identical (repeated structure)
  - <span style="color:red">ambiguous matches</span> are almost inevitable

- **Local search → sub-optimal solution**
  - take account of similarities between 2D-3D matches
  - involves in <span style="color:red">global compatibility</span> among all matching pairs? effective?

# Local search V.S. Global search

# How? – Step 1 (Build a map graph)

☐ Traditionally, 3D map is in the form of unordered point clouds

☐ Transform into weighted and bi-directed map graph

  ☐ covisibility $c_{ij}$ – if $P_j$ is seen by an image set, how likely $P_i$ can also be seen from the same image set

$$c_{ij} = \frac{|A_i \cap A_j|}{|A_j|}$$

# How? – Step 1 (Build a map graph)

☐ Collect all $c_{ij}$ into a square matrix C = $[c_{ij}]$ of size $N \times N$

☐ Normalize each column unit norm to form a left stochastic matrix

   ☐ with each column summing to 1
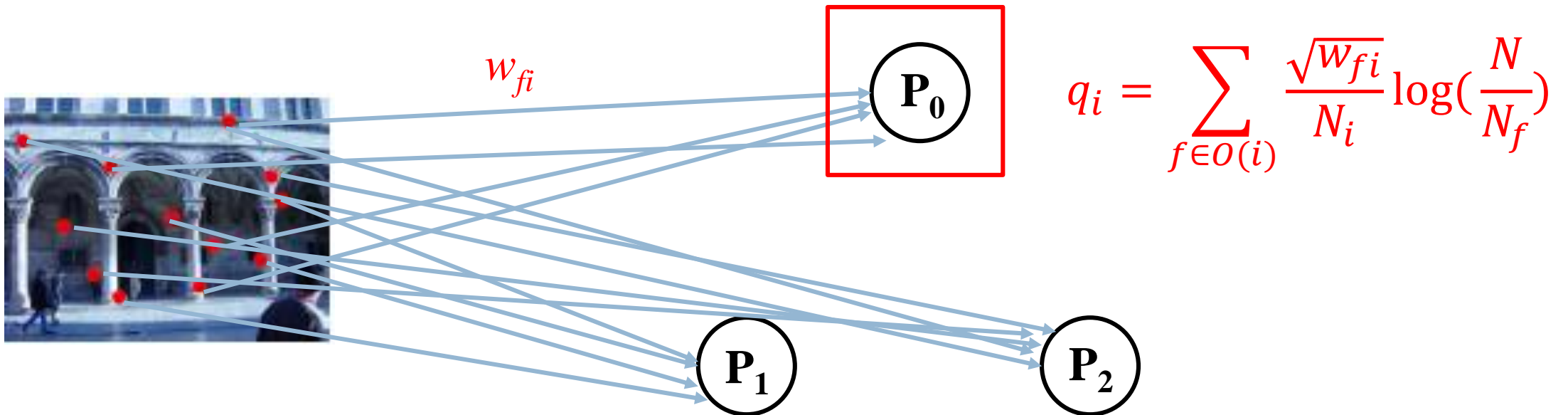
☐ Also call $C$ as state transition matrix

$$c_{ij} = \frac{\left|A_i \cap A_j\right|}{\left|A_j\right|} \quad \Rightarrow \quad \begin{bmatrix} c_{00} & \cdots & c_{0N} \\ \vdots & \ddots & \vdots \\ c_{N0} & \cdots & c_{NN} \end{bmatrix}$$

# How? – Step 2 (query vector)

□ query image → a set of 2D feature points

□ for every 2D feature points, find a set of tentative matches from 3D graph nodes, by comparing their descriptor similarity via Bag-of-words vocabulary tree search



$$q_i = \sum_{f \in O(i)} \frac{\sqrt{w_{fi}}}{N_i} \log(\frac{N}{N_f})$$

# How? – Step 2 (query vector)

- Collect all $q_i$ into a vector $q$
- Normalize $q$ to have unit norm $\quad q_i \leftarrow \dfrac{q_i}{\sum_{i=1}^{N} q_i}$

- $q$ can also be interpreted as a probability
  - measures the probability of point $i$ belongs to the optimal sub-set of 3D points that can be matched to the set of 2D query features

# How? – Step 3 (random walk)

- **Map graph → Markov network (aka. Markov Random Field)**
  - for global match between 2D query image and 3D map
  - when random walks converge, reach steady state
  - $p_v(t)$ is the probability of finding random walker at node $v$ at time $t$
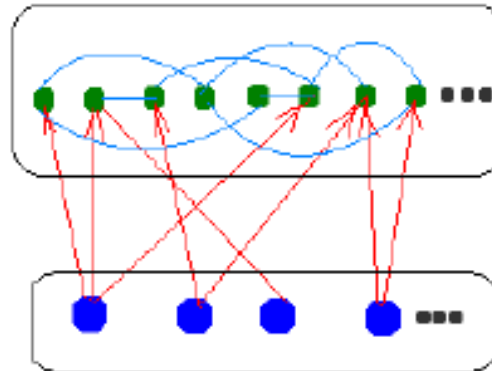  - $p_v(\infty)$ gives the probability that the random walker eventually ends at node $v$

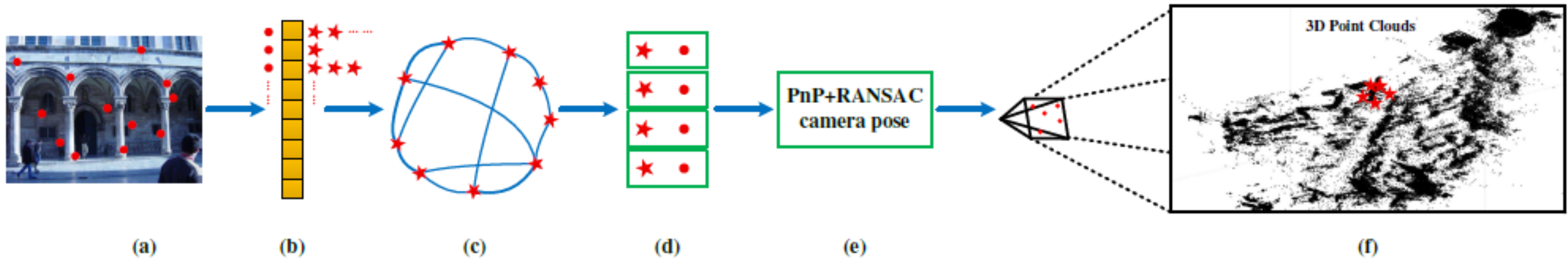# How? – Step 3 (random walk)

□ **Random Walk with Restart (RWR)**

▫ $p(t+1) = \alpha C p(t) + (1-\alpha)q$

▫ $p(0)$ is initialized by $q$

▫ $p(\infty)$ stands for the final matchability of every 3D point to the set of 2D query features

▫ sort $p(\infty)$ in descending order

# How? – Step 4 (camera pose computation)

☐ **Recover one-to-one correspondences**

  ☐ do <span style="color:red">ratio test</span> to retrieve one-to-one matches

  ☐ fed into <span style="color:red">PnP-RANSAC</span> to find camera position and orientation

# How? – Summary



(a)     (b)     (c)     (d)     (e)     (f)
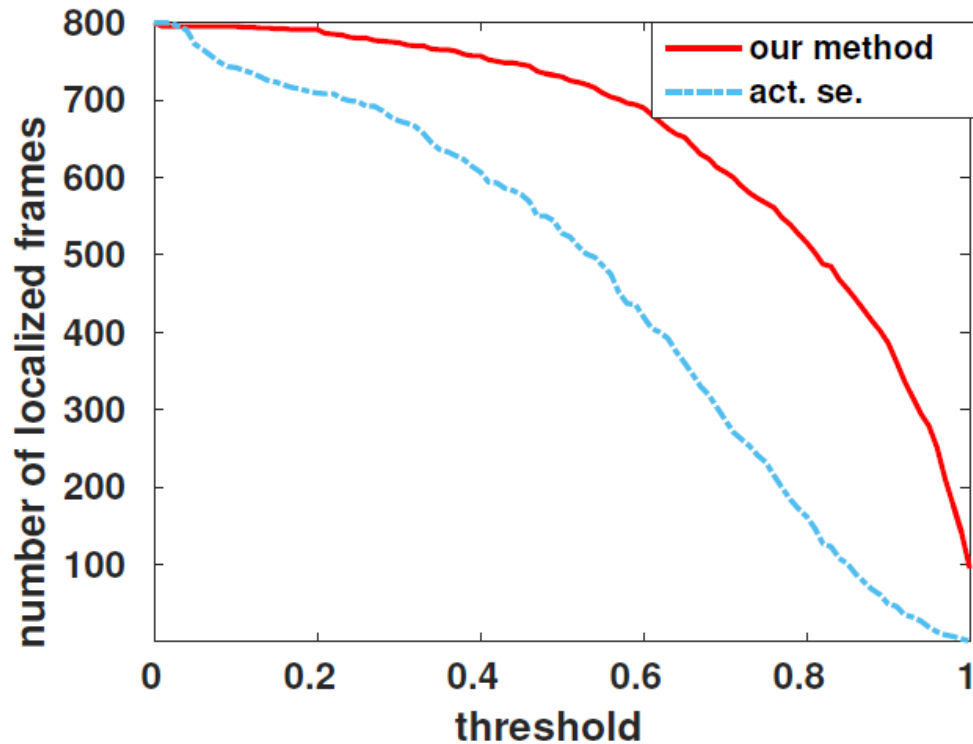
PnP+RANSAC camera pose

3D Point Clouds

# Experiment – dataset selection

☐ Four standard publicly available benchmark datasets for city-scale localization
(1) Dubrovnik, (2) Rome, (3) Vienna, (4) San Francisco (SF-0)

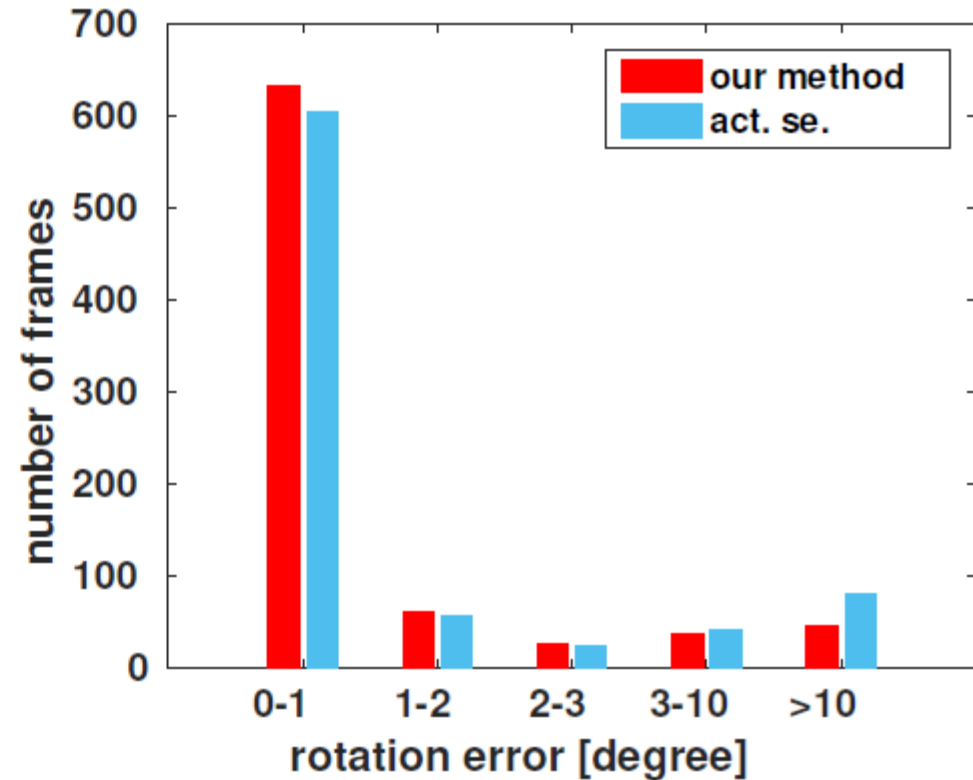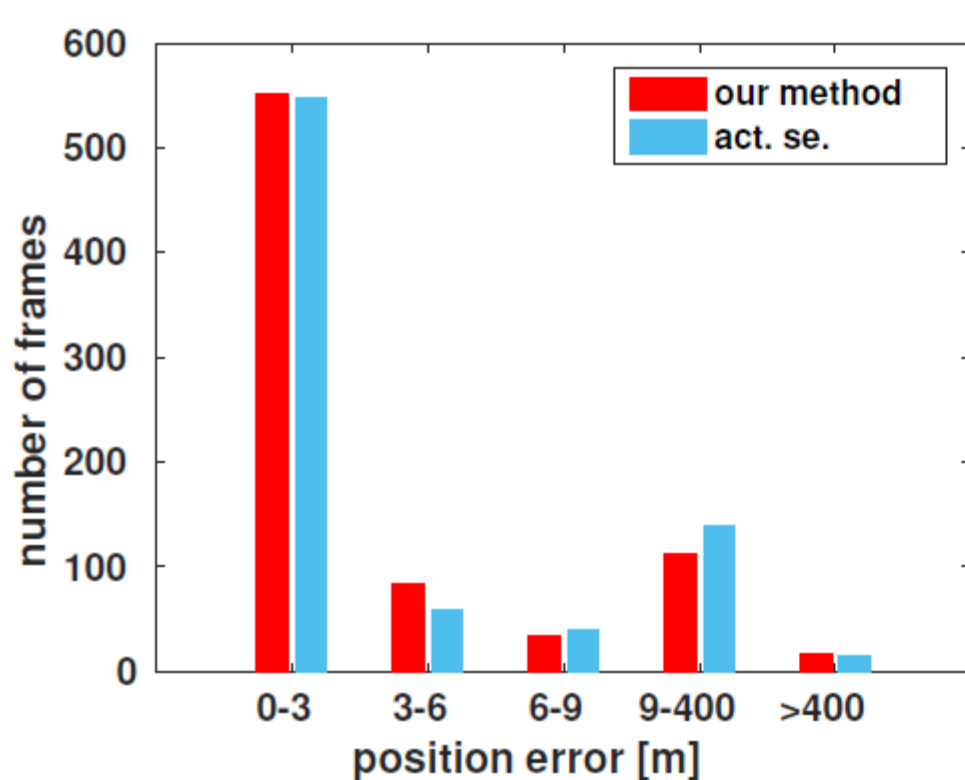| Dataset | #(images) | #(points) | #(query images) |
|---|---|---|---|
| Dubrovnik [32] | 6,044 | 1,975,263 | 800 |
| Rome [32] | 15,179 | 4,067,119 | 1,000 |
| Vienna [23] | 1,324 | 1,123,028 | 266 |
| SF-0 [12] | 610,773 | 30,342,328 | 803 |

# Experiment

- In term of recall-rate (# images have been successfully localized)



| Method | Inlier thresholds | | | | | | |
|---|---|---|---|---|---|---|---|
| | 0.2 | 0.3 | 0.4 | 40.5 | 0.6 | 0.7 | 0.8 |
| Active Search [43, 44] | 709 | 673 | 607 | 528 | 420 | 287 | 162 |
| **Our method** | **791** | **774** | **757** | **730** | **690** | **607** | **516** |

# Experiment
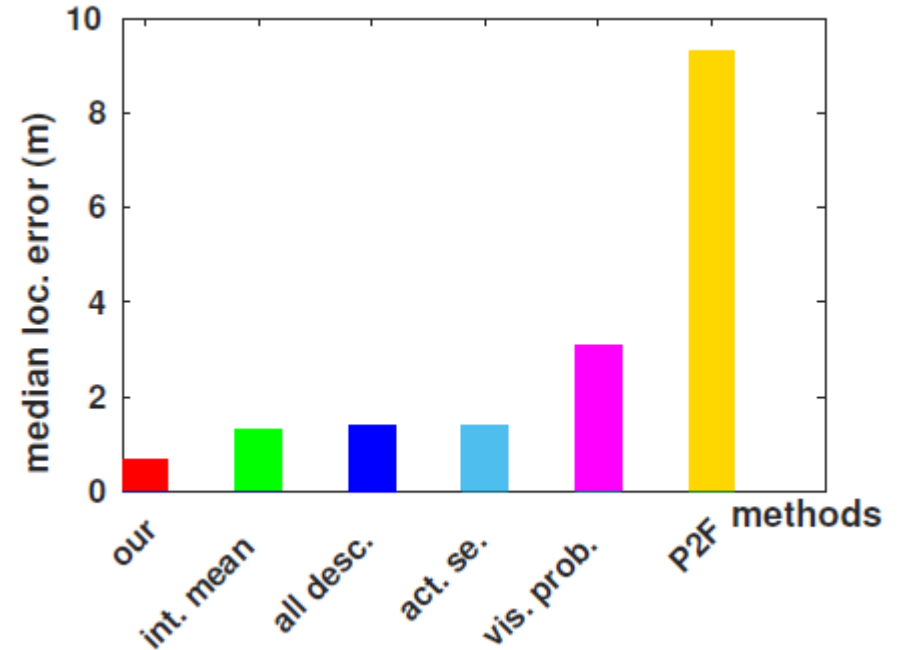
☐ In term of precision (camera localization error)

# Experiment

- In term of precision (camera localization error)

| Method | quartile errors (m) | | | num. of images | | |
|---|---|---|---|---|---|---|
| | 1st | median | 3rd | <18.3m | >400m | #(reg.) |
| **our method** | **0.24** | **0.70** | **2.67** | **743** | **7** | 794 |
| act. se. [43,44] | 0.40 | 1.40 | 5.30 | 704 | 9 | **795** |
| all desc. [42] | 0.40 | 1.40 | 5.90 | 685 | 16 | 783 |
| int. mean [42] | 0.50 | 1.30 | 5.10 | 675 | 13 | 782 |
| P2F [32] | 7.50 | 9.30 | 13.40 | 655 | - | 753 |
| vis. prob. [13] | 0.88 | 3.10 | 11.83 | - | - | 788 |

# Experiment

- Demo

# Conclusion

- a global method, taking account of not only individual feature match's visual similarity but also the global compatibilities as measured by the pair-wise covisibility, to deal with scalability and ambiguity for localization