# MatchMiner: Efficient Spanning Structure Mining in Large Image Collections
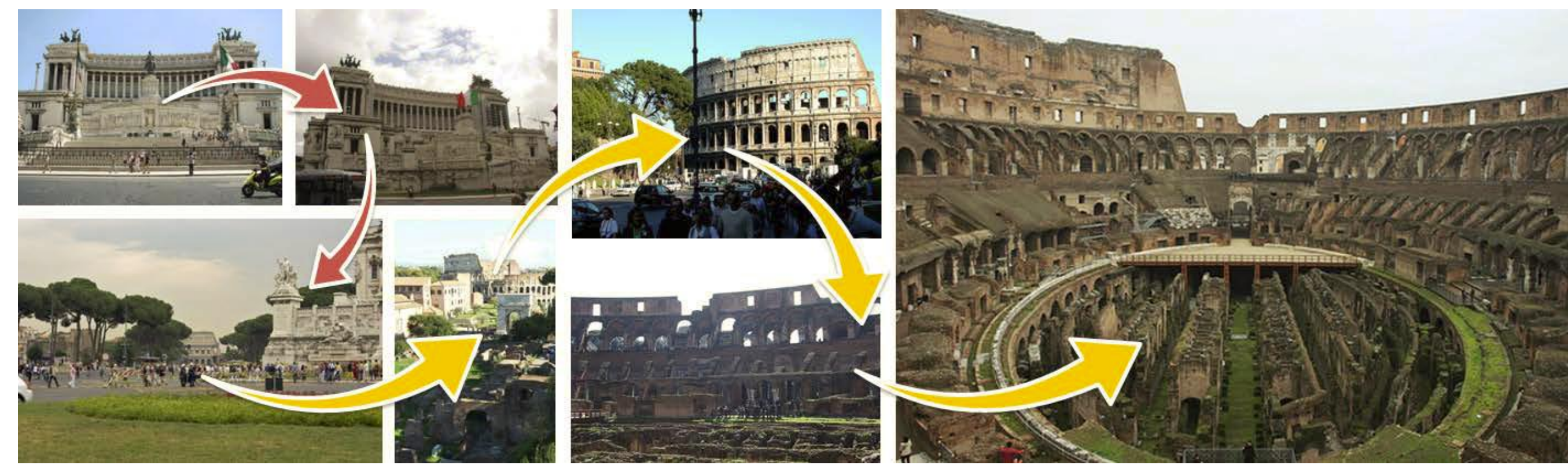
Yin Lou, Noah Snavely, and Johannes Gehrke
Department of Computer Science, Cornell University

12th European Conference on Computer Vision
Firenze, Italy, 7-13 October 2012

## Motivation

- Internet photos cover large parts of the world
- Novel applications are using **image graphs**
- We want to **connect** images as efficiently as possible
- We focus on finding connected components



## Challenges with Unstructured Collections

- Image matching is **expensive**
- It is hard to know promising image pairs beforehand
- Visual similarity is a noisy predictor
- Large image collections have many "singleton" images

## Contributions: a large-scale image matcher that:

- We incorporate relevance feedback
- We propose rank distance to prune singleton images
- We propose an information-theoretic approach

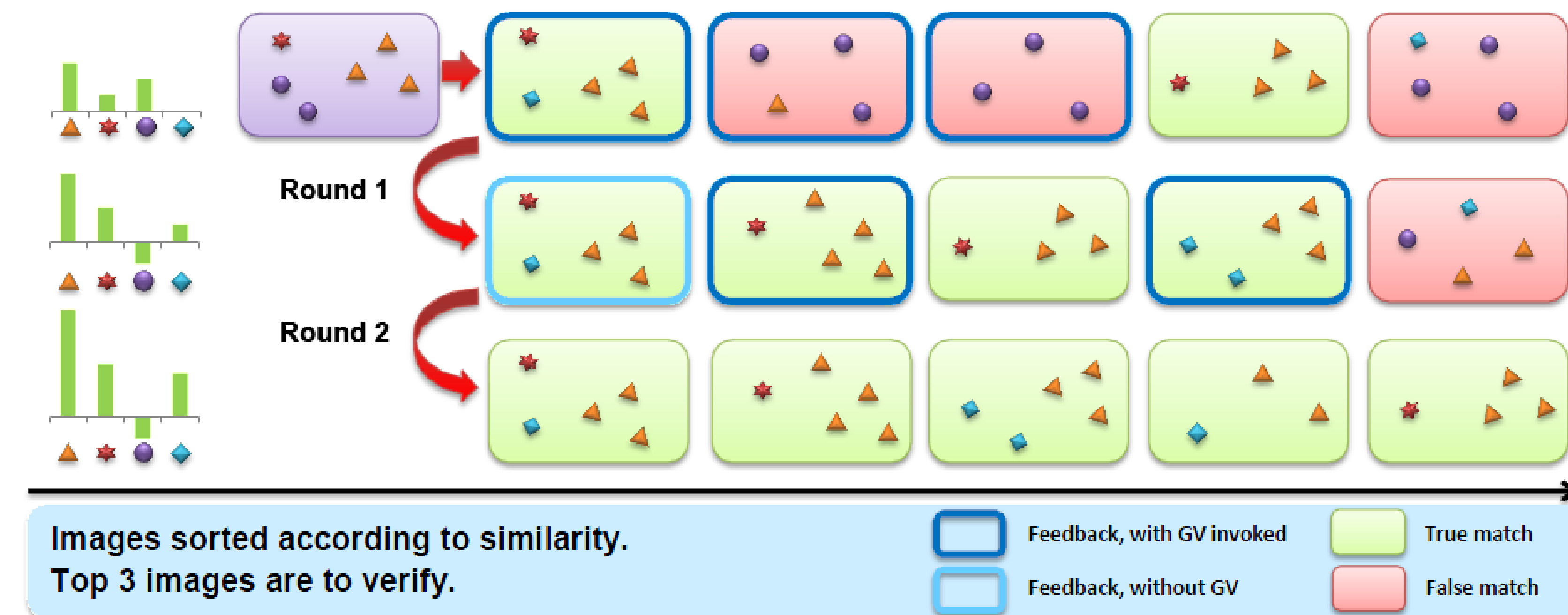## Image Representation and Matching Procedure

- Each image is represented using BoW model
- One million visual words are trained offline
- Standard tf-idf weights are applied on image vectors
- We use standard geometric verification procedure
  - SIFT matching
  - RANSAC-based F-matrix estimation

## MatchMiner

Two stage approach: (1) we find an initial set of CCs by matching similar images, incorporating relevance feedback, (2) we merge CCs using an information-theoretic approach and discard singleton images.
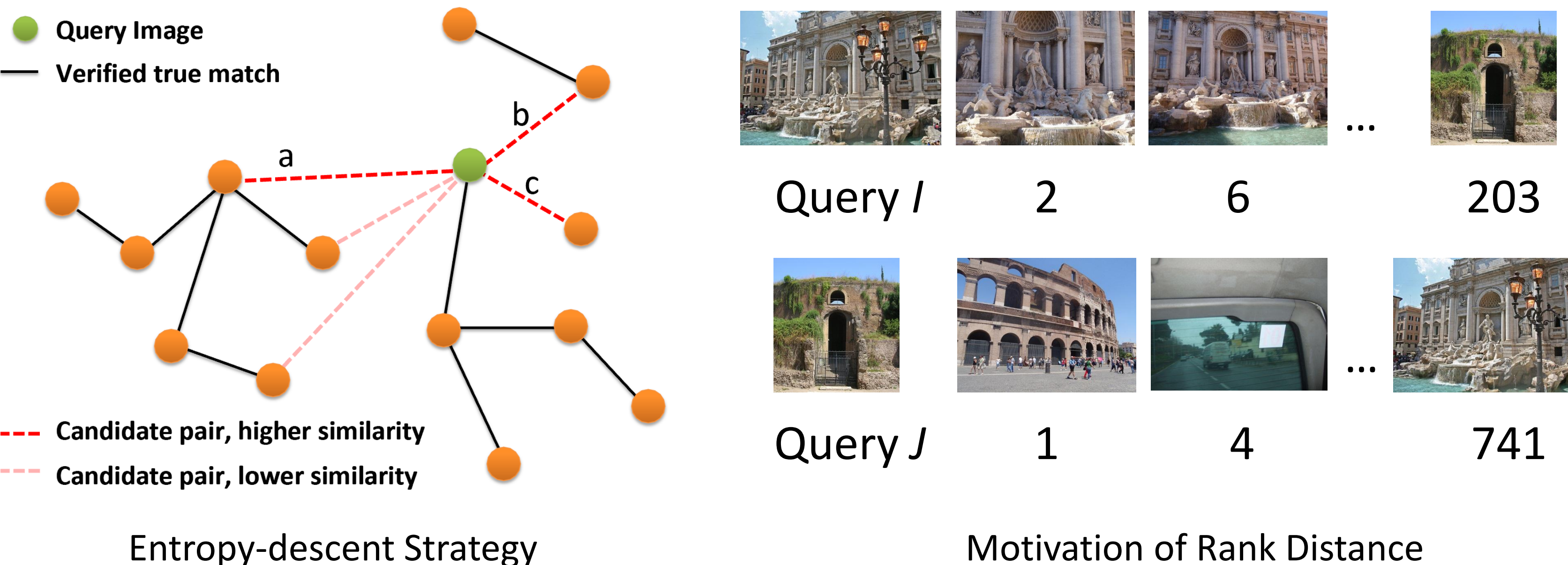
### Step 1

- Each image vector $Q_1$ retrieves a short list of images $\{I\}$
- Geometric verification partitions $\{I\}$ into two sets, $P$ and $N$
- Relevance feedback: $Q_{t+1} = Q_t + \alpha^{t+1}/|P| \sum_{I \in P} I - \beta^{t+1}/|N| \sum_{I \in N} I$



Round 1
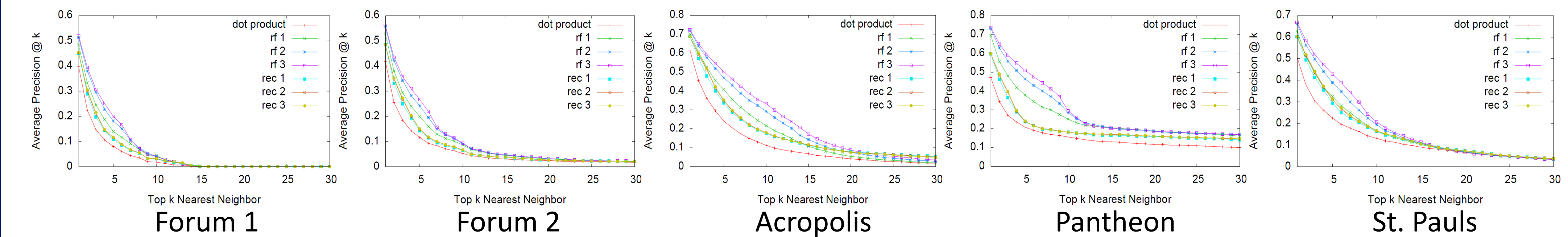Round 2

Images sorted according to similarity. Top 3 images are to verify.

Feedback, with GV invoked
Feedback, without GV
True match
False match

### Step 2

- Minimizing entropy H(C); prefer to merge large CCs
- Rank distance: $R(I, J) = 2 Rank_I(J) Rank_J(I)/(Rank_I(J) + Rank_J(I))$



Query $I$    2    6    ...    203

Query $J$    1    4    ...    741

- - - Candidate pair, higher similarity
- - - Candidate pair, lower similarity

Query Image
Verified true match

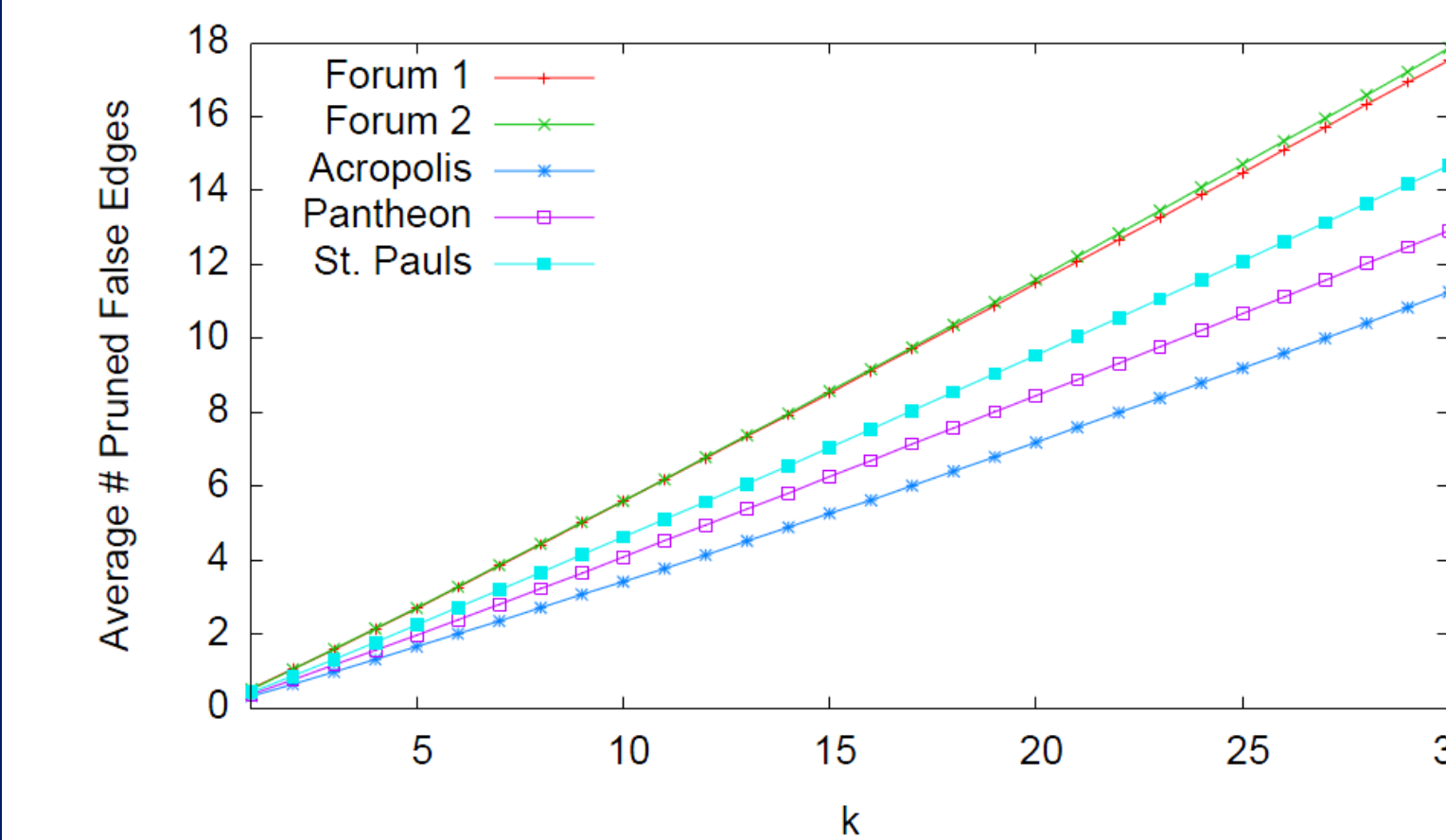Entropy-descent Strategy | Motivation of Rank Distance

## Experiments

- Five medium-sized datasets and two large datasets
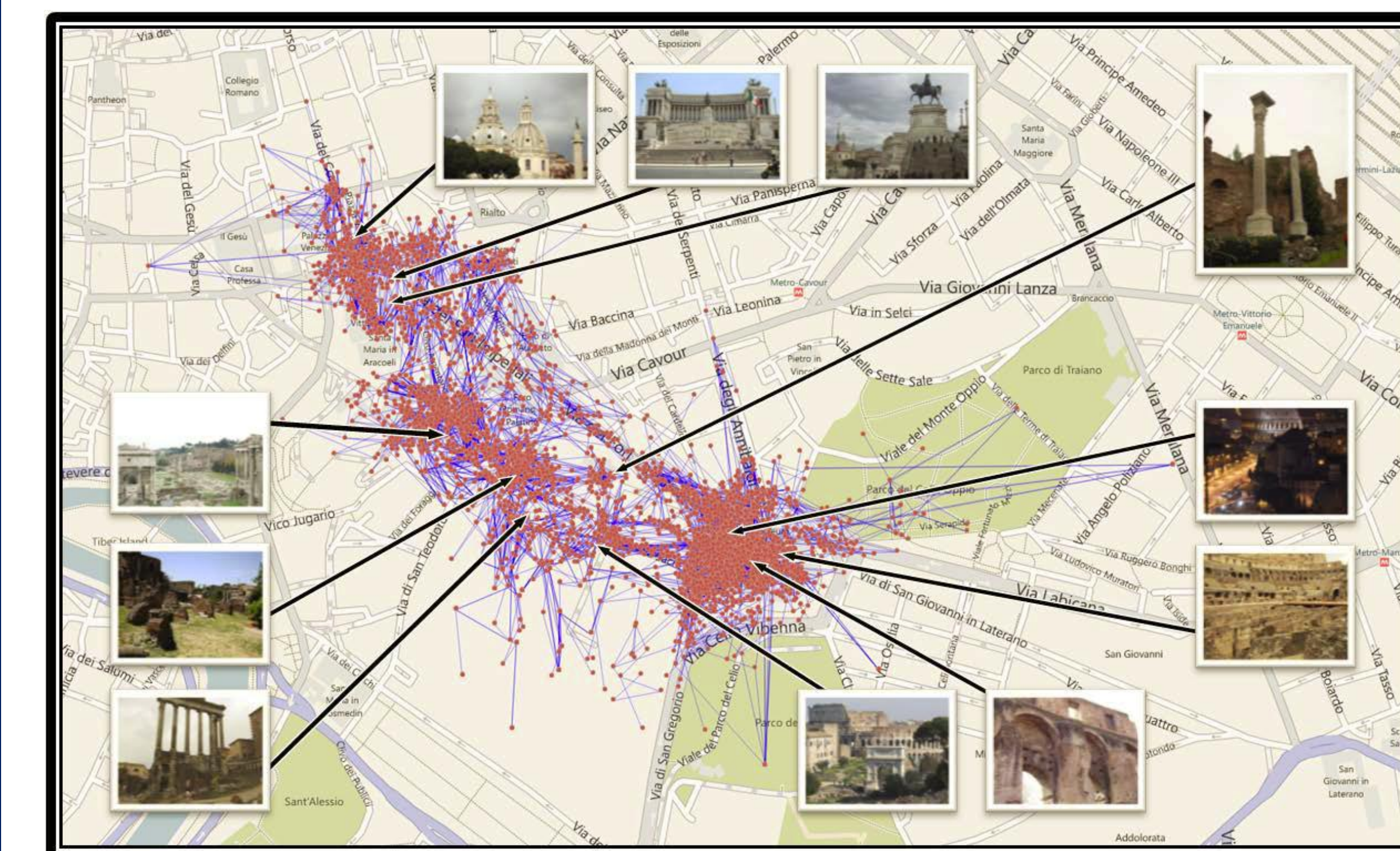- We compare MatchMiner with Image Webs [Heath et al. 10]

### Relevance Feedback



Forum 1    Forum 2    Acropolis    Pantheon    St. Pauls

### Rank Distance



False Edges Pruned by RD
Rate of prunning true edges <0.1%

### Mining Results

| Dataset | Algorithm | $CC_1$ | $\overline{MI}$ |
|---------|-----------|--------|-----|
| Forum1 | ImageWebs | 105 | 0.80 |
| | MatchMiner | 266 | 0.90 |
| Forum 2 | ImageWebs | 728 | 0.80 |
| | MatchMiner | 908 | 0.85 |
| Acropolis | ImageWebs | 1894 | 0.79 |
| | MatchMiner | 1948 | 0.83 |
| Pantheon | ImageWebs | 639 | 0.59 |
| | MatchMiner | 765 | 0.74 |
| St. Pauls | ImageWebs | 1816 | 0.80 |
| | MatchMiner | 1883 | 0.84 |

K = 20

| Dataset | Algorithm | $CC_1$ | $\overline{MI}$ |
|---------|-----------|--------|-----|
| Forum1 | ImageWebs | 180 | 0.85 |
| | MatchMiner | 271 | 0.91 |
| Forum 2 | ImageWebs | 788 | 0.83 |
| | MatchMiner | 937 | 0.87 |
| Acropolis | ImageWebs | 1951 | 0.84 |
| | MatchMiner | 1978 | 0.86 |
| Pantheon | ImageWebs | 659 | 0.62 |
| | MatchMiner | 788 | 0.80 |
| St. Pauls | ImageWebs | 1845 | 0.82 |
| | MatchMiner | 1934 | 0.89 |

K = 30

### Mining Large-scale Datasets



- Largest CC of Forum
- 1 hr 39 min
- 53 nodes

| Dataset | Algorithm | $CC_1$ | $H(C)$ |
|---------|-----------|--------|--------|
| Forum | ImageWebs | 6944 | 11.92 |
| | MatchMiner | 13871 | 11.62 |
| Washington DC | ImageWebs | 11249 | 16.76 |
| | MatchMiner | 16922 | 16.64 |