

Cross-View Geo-Localization: Ground-to-Aerial Image Matching

Mubarak Shah

shah@crcv.ucf.edu

Center for Research in Computer Vision

University of Central Florida

<http://www.crcv.ucf.edu>

Geo-Localization

- Pixel-Wise Geo-Localization
 - Given a **query** image, geo-localize each **pixel** by **aligning** an image with the geodetically accurate **reference** image.
- Image-Based Geo-Localization
 - Given a **query** image, determine its GPS location by **matching** it with **geo-tagged** reference images.

Contents

- Pixel-Wise Geo-localization
 - Geodetic Alignment of Aerial Video Frames
- Image-Based Geo-Localization
 - Same View (Street-View to Street-View)
 - Generalized Maximum Clique (PAMI, 2014)
 - Constraint Dominant Sets (PAMI, 2017)
 - Cross-View Geo-Localization
 - Bird's Eye-View to Street View (CVPR, 2017)
 - Aerial to Ground View (ICCV, 2019)

Geodetic Alignment of Aerial Video Frames

Y. Sheikh, S. Khan, M. Shah and R. Cannata
2003



Associate Professor
CMU/FaceBook

*geo*Registration

Data Overview

Aerial Video Data

Reference Data

Telemetry



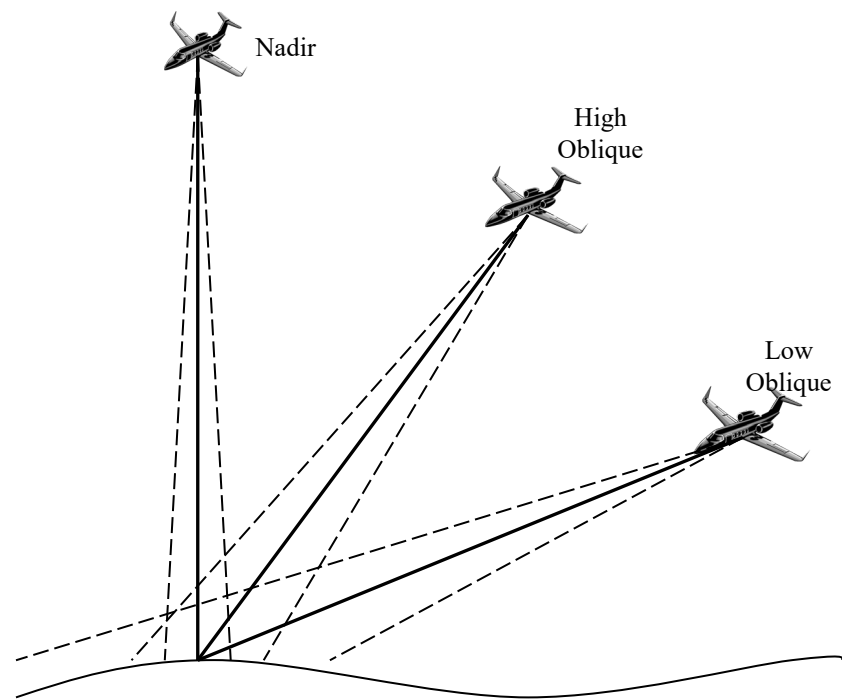
*geo*Registration

Different Viewpoints

Nadir

High Oblique

Low Oblique



*geo*Registration

Aerial Video Imagery

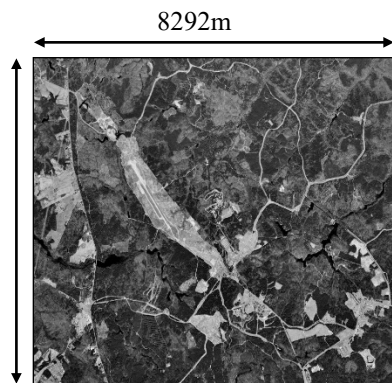


geoRegistration

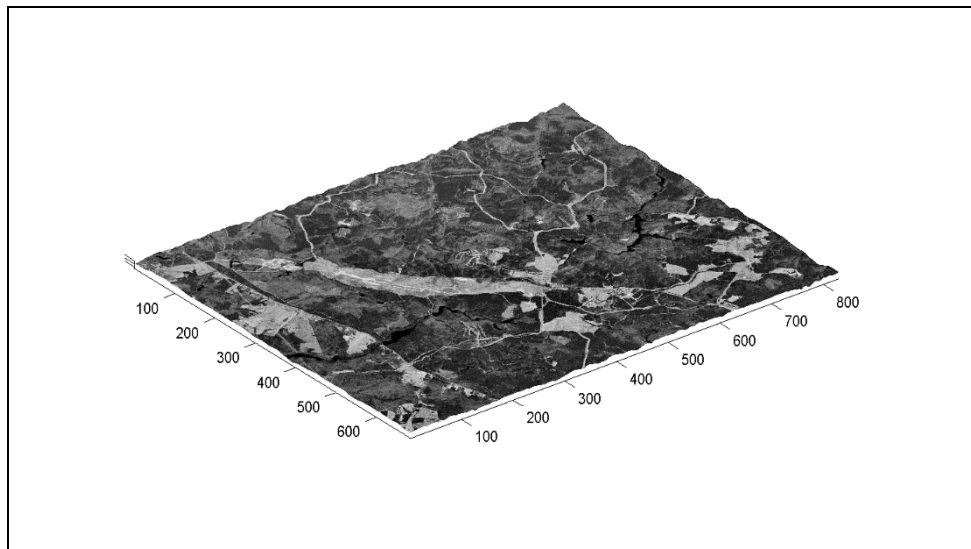
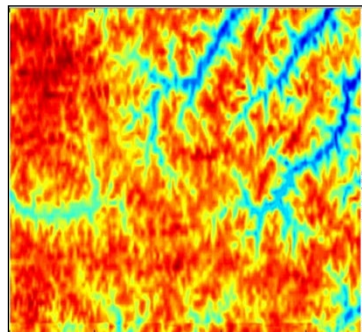
Reference Data

DOQ
(Digital Ortho Quad)

6856 m



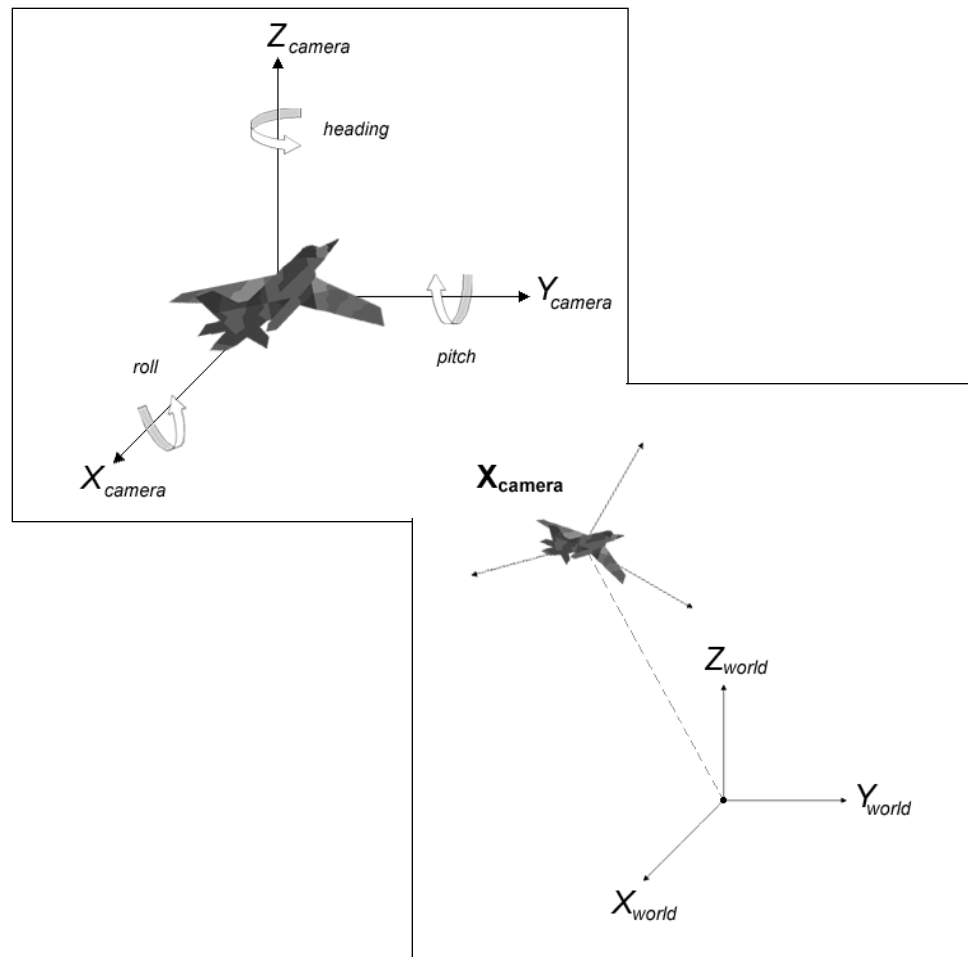
DEM
(Digital Elevation Map)



geoRegistration

Telemetry Data

- { Vehicle Longitude
- { Vehicle Latitude
- { Vehicle Height
- { Vehicle Heading
- { Vehicle Roll
- { Vehicle Pitch
- { Camera Elevation Angle
- { Camera Scan Angle
- { Camera Focal Length

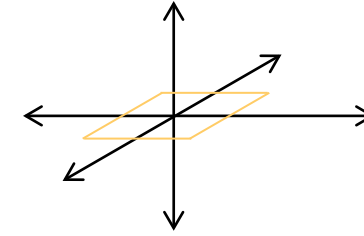


Sensor Model

$$\mathbf{X}_{\text{camera}} = [X_{\text{camera}} \quad Y_{\text{camera}} \quad Z_{\text{camera}}]^T$$

$$\mathbf{X}_{\text{world}} = [X_{\text{world}} \quad Y_{\text{world}} \quad Z_{\text{world}}]^T$$

$$\mathbf{X}_{\text{camera}} = \Pi_t \mathbf{X}_{\text{world}}$$



$\Pi_t = f$ (camera elevation, camera scan, vehicle pitch, vehicle roll, vehicle heading, vehicle elevation)

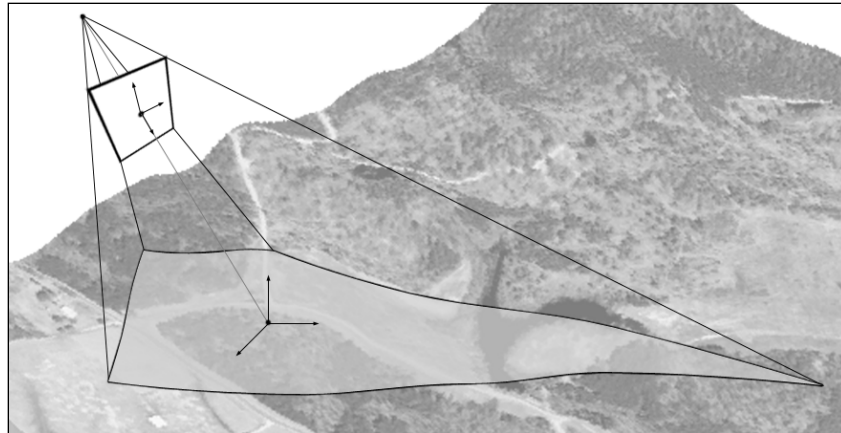
$$\Pi_t = \begin{bmatrix} \cos \omega & 0 & -\sin \omega & 0 \\ 0 & 1 & 0 & 0 \\ \sin \omega & 0 & \cos \omega & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \tau & \sin \tau & 0 & 0 \\ -\sin \tau & \cos \tau & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \phi & 0 & -\sin \phi & 0 \\ 0 & 1 & 0 & 0 \\ \sin \phi & 0 & \cos \phi & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos \beta & \sin \beta & 0 \\ 0 & -\sin \beta & \cos \beta & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \alpha & \sin \alpha & 0 & 0 \\ -\sin \alpha & \cos \alpha & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & \Delta T_x \\ 0 & 1 & 0 & \Delta T_y \\ 0 & 0 & 1 & \Delta T_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Sensor Model

$$\Pi_c = P\Pi_t$$

$$\Pi_c = PG_yG_zR_yR_xR_zT$$

$$\mathbf{x}_{\text{video}} = \Pi_c \mathbf{X}_{\text{world}}$$



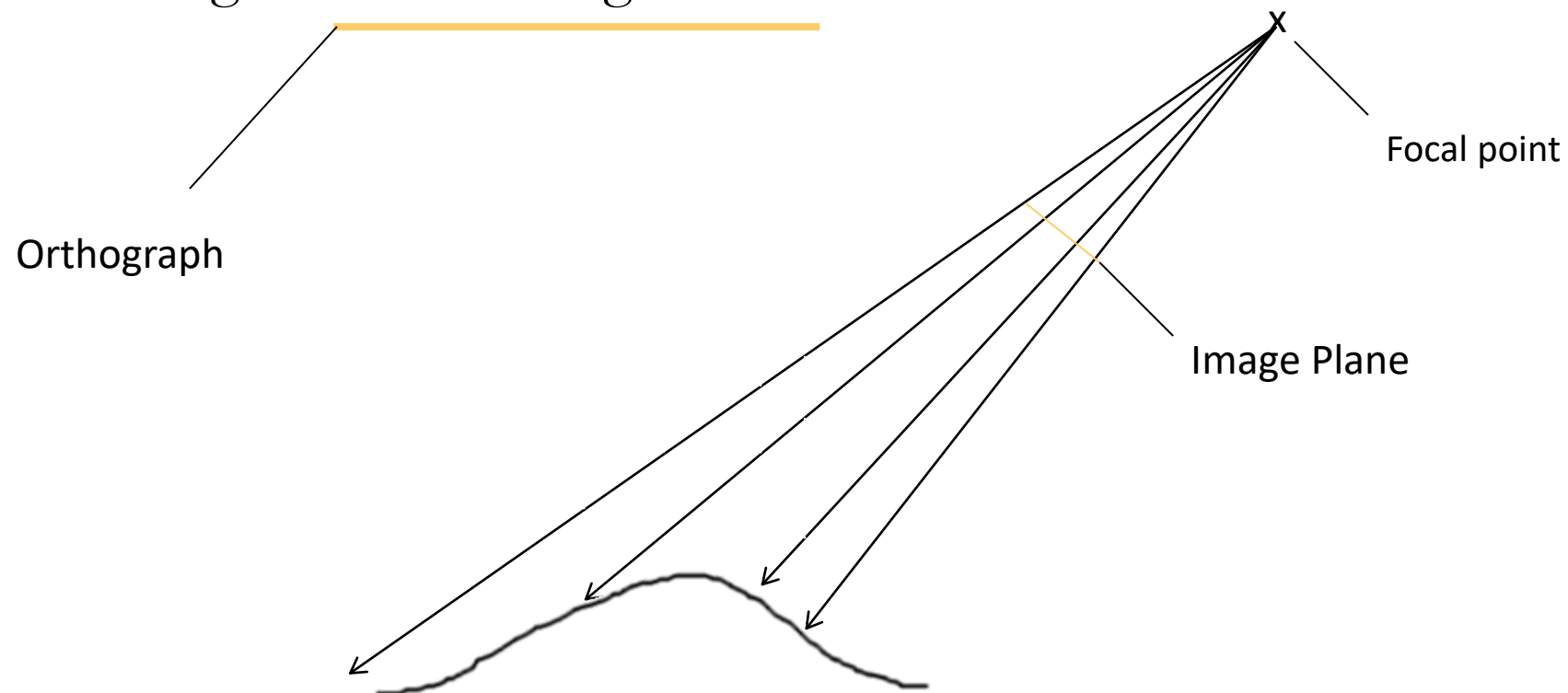
*ortho*Rectification

- Bringing both imageries onto a common view projection (**Bridging the gap**)
- Accounts for gross misalignment
- Not accurate due to telemetry noise



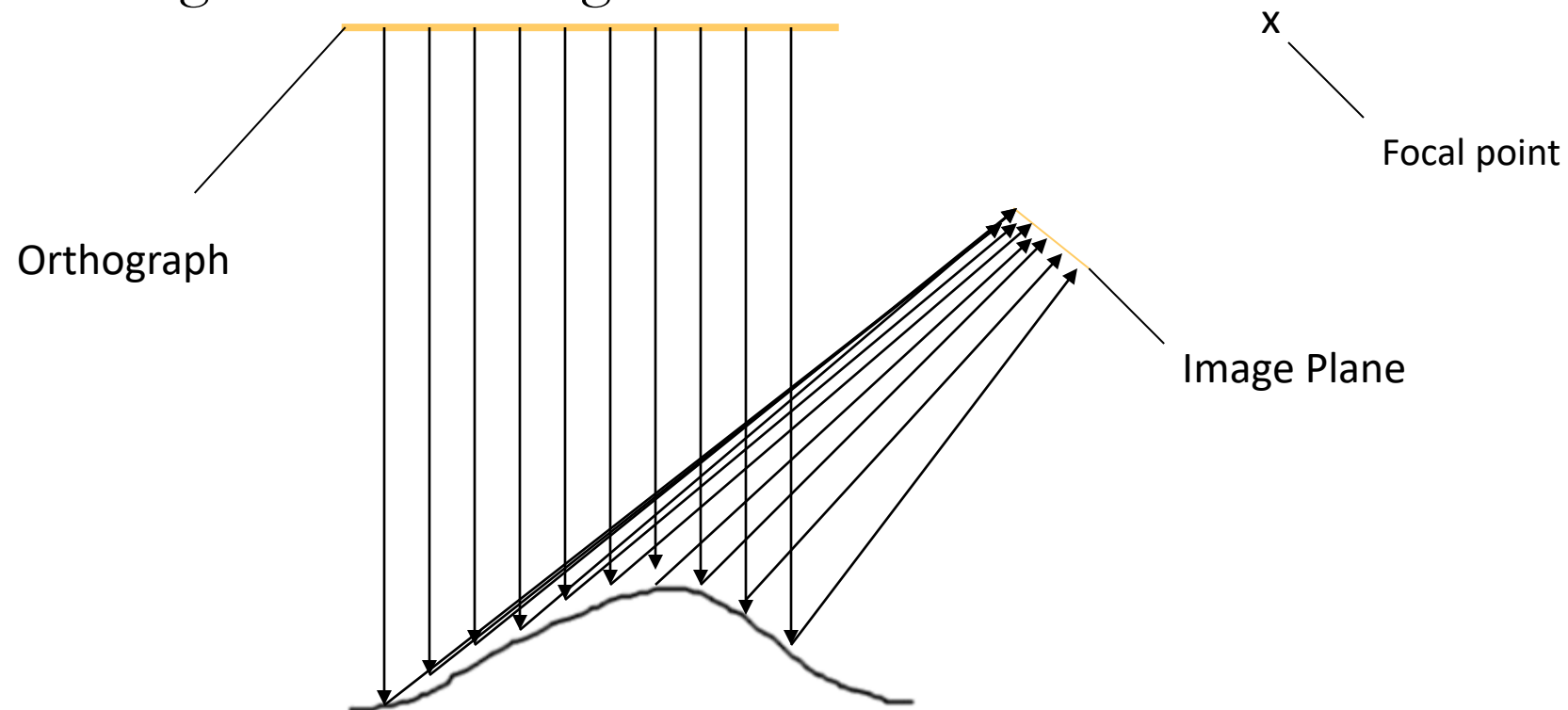
Rectification

Creating the OrthoImage



*ortho*Rectification

Creating the OrthoImage



orthoRectification

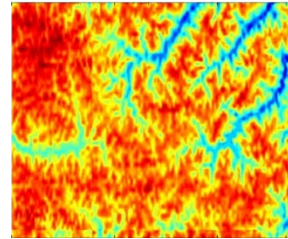
Sensor File + Mission Images + Reference Eniv. = Orthorectification

```
OTTER
system_id      TV
               sensor_type
               serial_number
0001
9.4000815266640300e+08  image_time
3.813183746469612200e+01  vehicle_latitude
-7.734523185193877700e+01  vehicle_longitude
9.9496284098763800e+02  vehicle_height
9.9951717444103990e+01  vehicle_pitch
1.70162641811320900e+00  vehicle_roll
1.20701055175302940e+02  vehicle_heading
1.658968732990974800e+02  camera_focal_length
-5.351314188957229100e+01  camera_elevation
-7.22969419154670500e+00  camera_scan_angle
480      number_image_lines
640      number_image_samples
1.0000000000000000e+00  v_x_scale_factor
1.0000000000000000e+00  v_y_scale_factor
3.812941913638164900e+01  ground_ref_gps_lat
-7.733710187470967400e+01  ground_ref_gps_lon
-1.4180415664242940e+09  ground_ref_gps_hgt
1.73205080756887200e+02  position_error
1.0000000000000000e+02  vehicle_lat_sigma
1.0000000000000000e+02  vehicle_lon_sigma
1.0000000000000000e+02  vehicle_hgt_sigma
5.0000000000000000e+01  vehicle_heading_sigma
5.0000000000000000e+01  vehicle_pitch_sigma
5.0000000000000000e+01  vehicle_roll_sigma
6.18638046322339200e+04  focal_length_sigma
5.0000000000000000e+01  gimbal_azimuth_sigma
5.0000000000000000e+01  gimbal_elevation_sigma
```

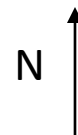
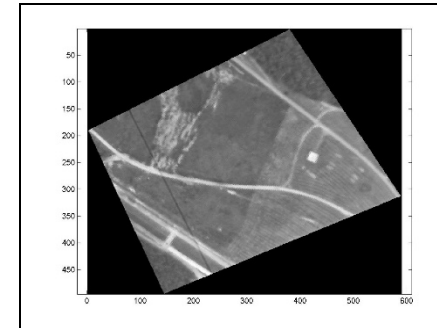
+



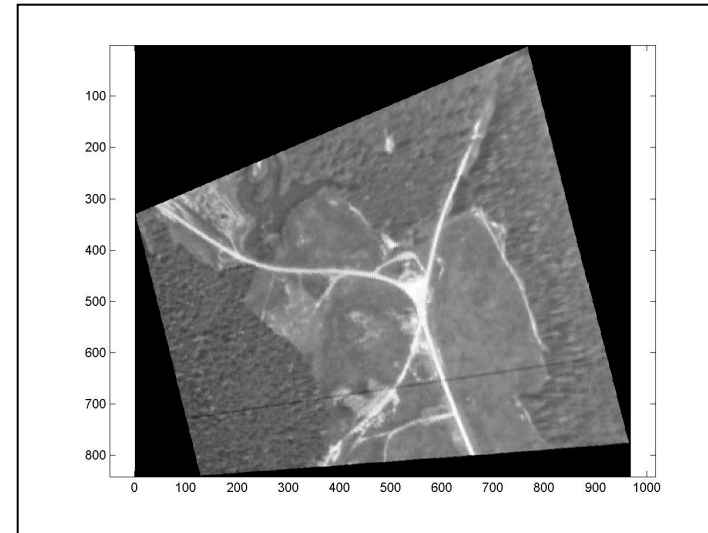
+



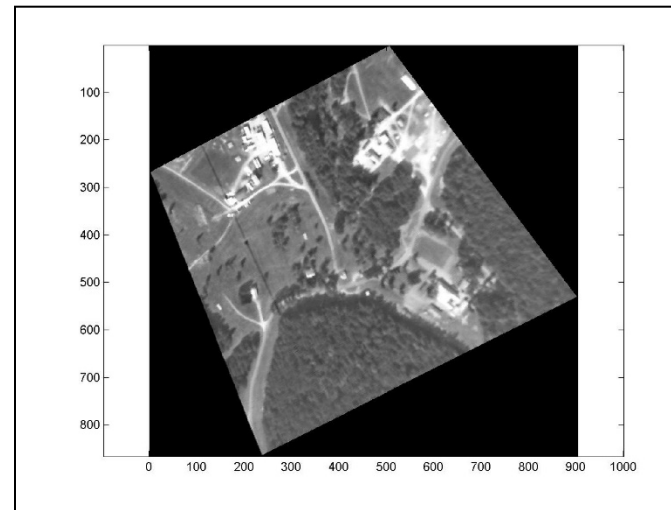
=



*ortho*Rectification



*ortho*Rectification



Reference Image
(DOQ)







Contents

- Pixel-Wise Geo-localization
 - Geodetic Alignment of Aerial Video Frames
- **Image-Based Geo-Localization**
 - Same View (Street-View to Street-View)
 - Generalized Maximum Clique (PAMI, 2014)
 - Constraint Dominant Sets (PAMI, 2017)
 - Cross-View Geo-Localization
 - Bird's Eye-View to Street View (CVPR, 2017)
 - Aerial to Ground View (ICCV, 2019)

Geo-localization Using Image Matching

- No Meta-Data (Telemetry)
- No Sensor Model
- No Geodetically Accurate Reference Image
 - Only Geo-tagged images
- Image level coarse geo-localization

“Where Am I?”

➤ Problem:

Image Localization

Input



Output



Mere Visual Information(Images) Location in Terms of λ (Lon.) and ϕ (Lat.)

“Where Am I?”

➤ Problem:

Image Localization

Input



Output



Mere Visual Information(Images) Location in Terms of λ (Lon.) and ϕ (Lat.)
 $\phi=40.4419, \lambda=-79.9986$

Google Maps Street View Dataset

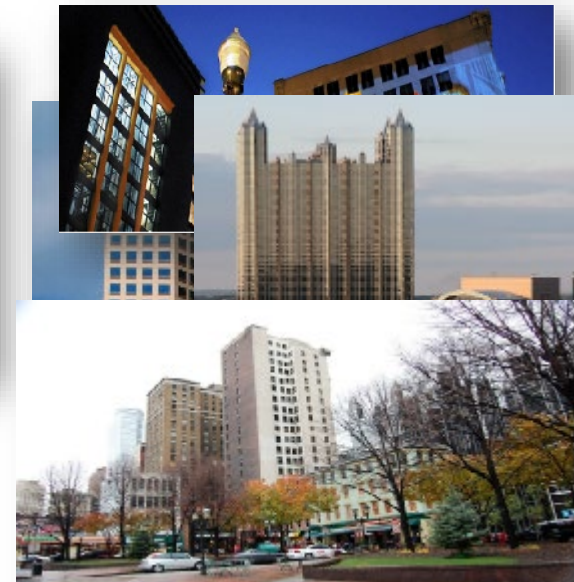
- Reference Set: GSV ~100k
- Test Set:
 - 521 GPS- Tagged from Pittsburgh, PA and Orlando, FL.
 - Downloaded From Flickr, Panoramio, Picasa, etc.



picasa.google.com



Panoramio.com

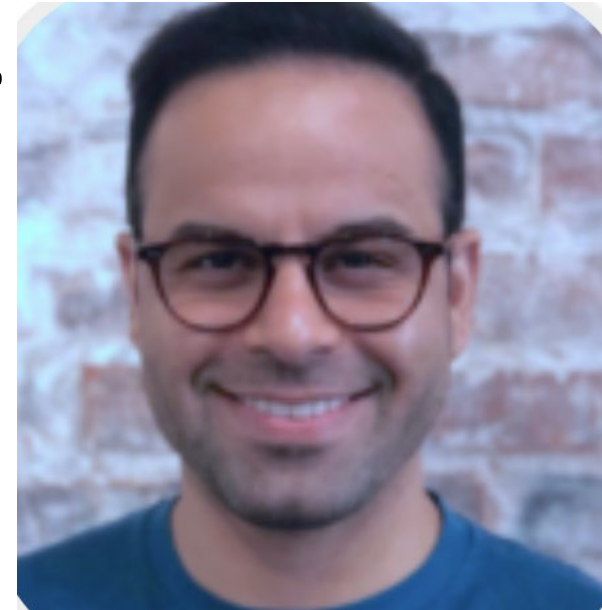


Flickr.com

Image Geo-localization Based on Multiple Nearest Neighbor Feature Matching Using Generalized Graphs.

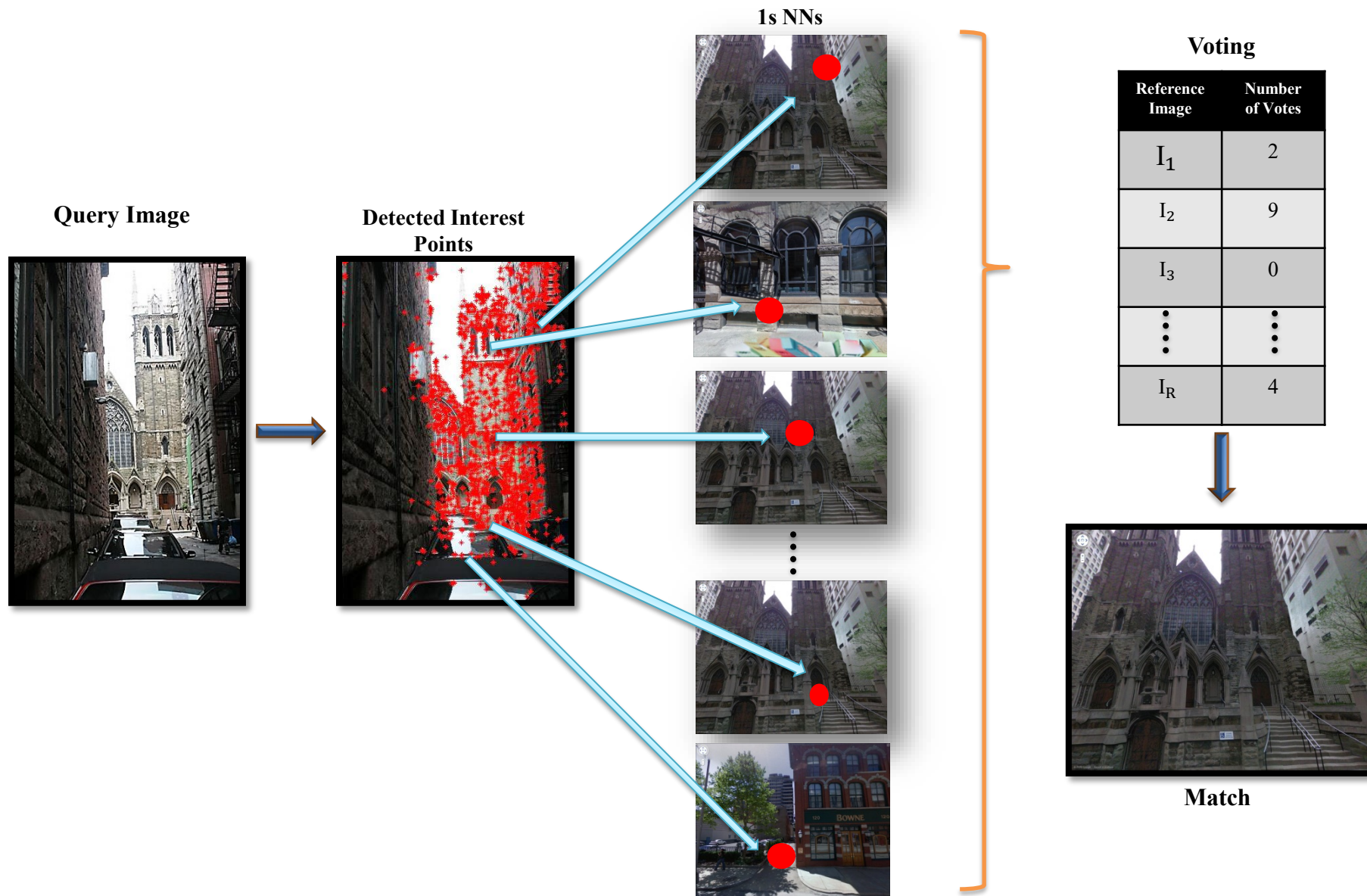
Amir Zamir and Mubarak Shah

In *T-PAMI*, 2014.

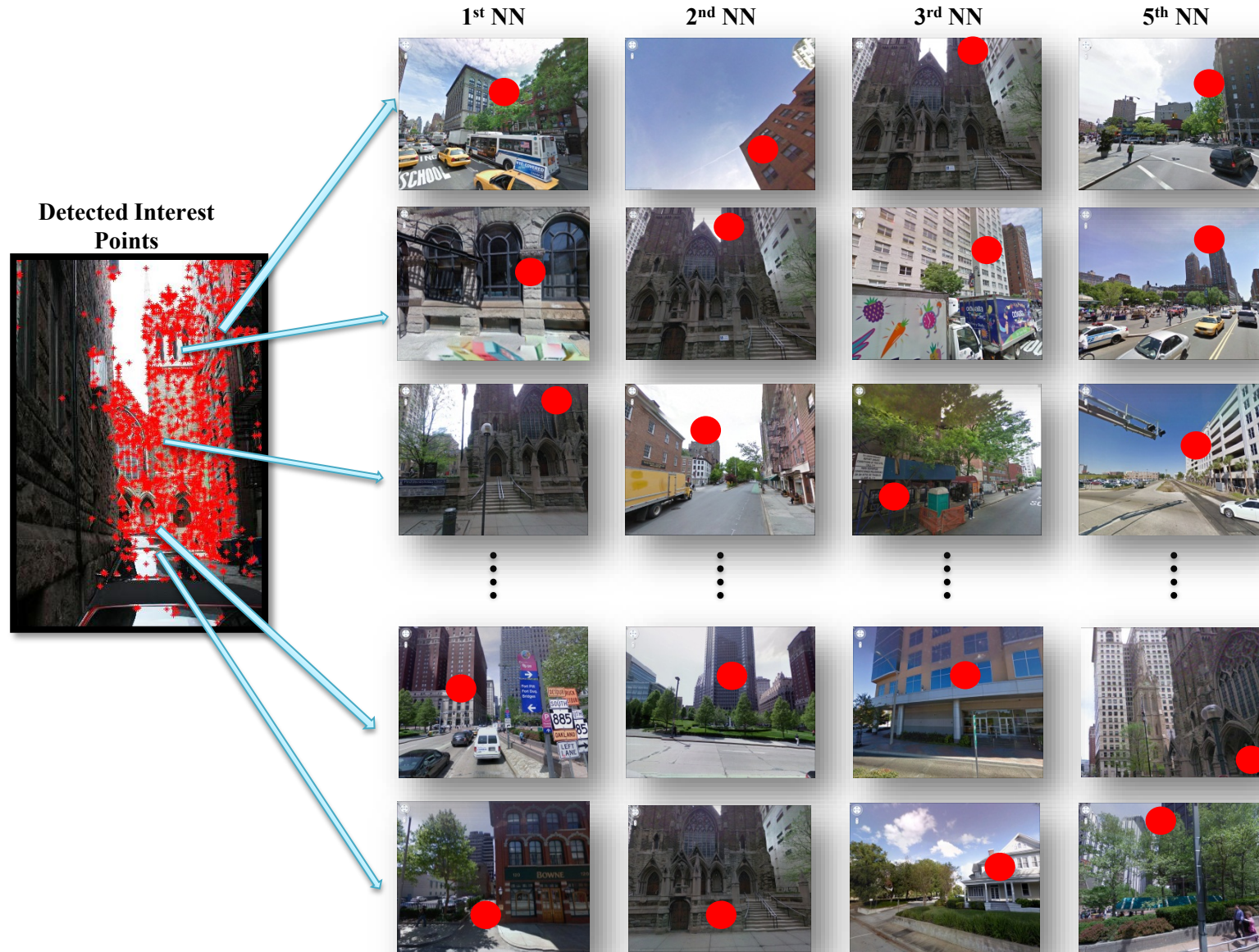


Asst Professor; EPFL

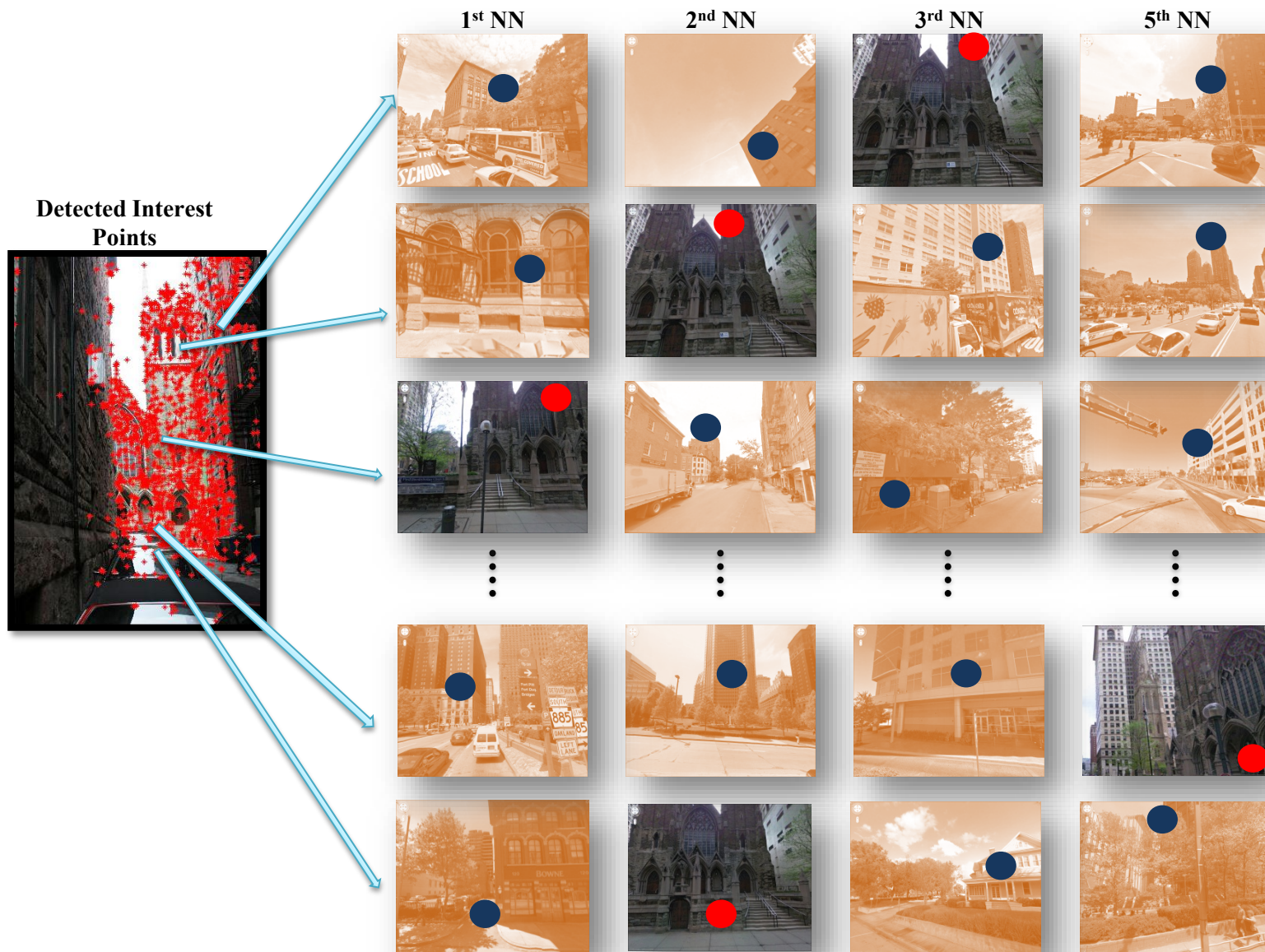
I: Localization Recognition Using Local Features



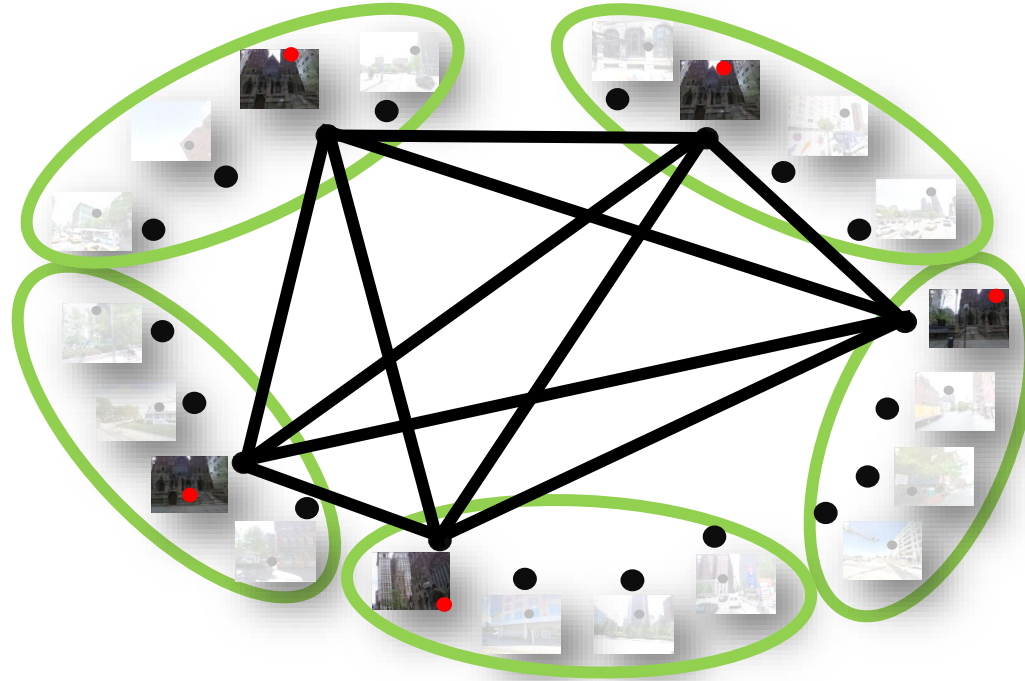
Using Multiple Nearest Neighbors



Using Multiple Nearest Neighbors



Generalized Minimum Clique



- Subset of NNs with maximum agreement in local and global features

Geo-localization Results



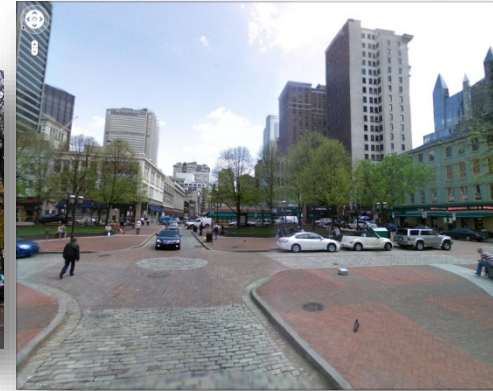
Query



Match – Error: 7.6 m



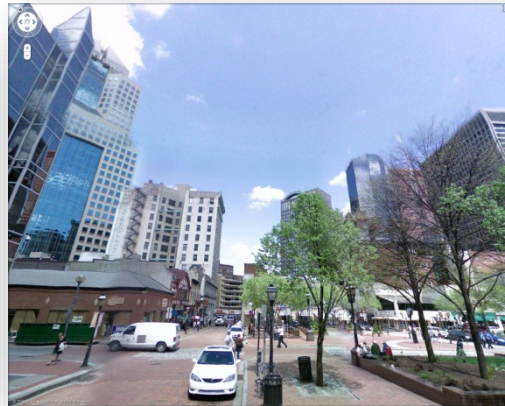
Query



Match – Error: 6.9 m



Query



Match – Error: 308.1 m

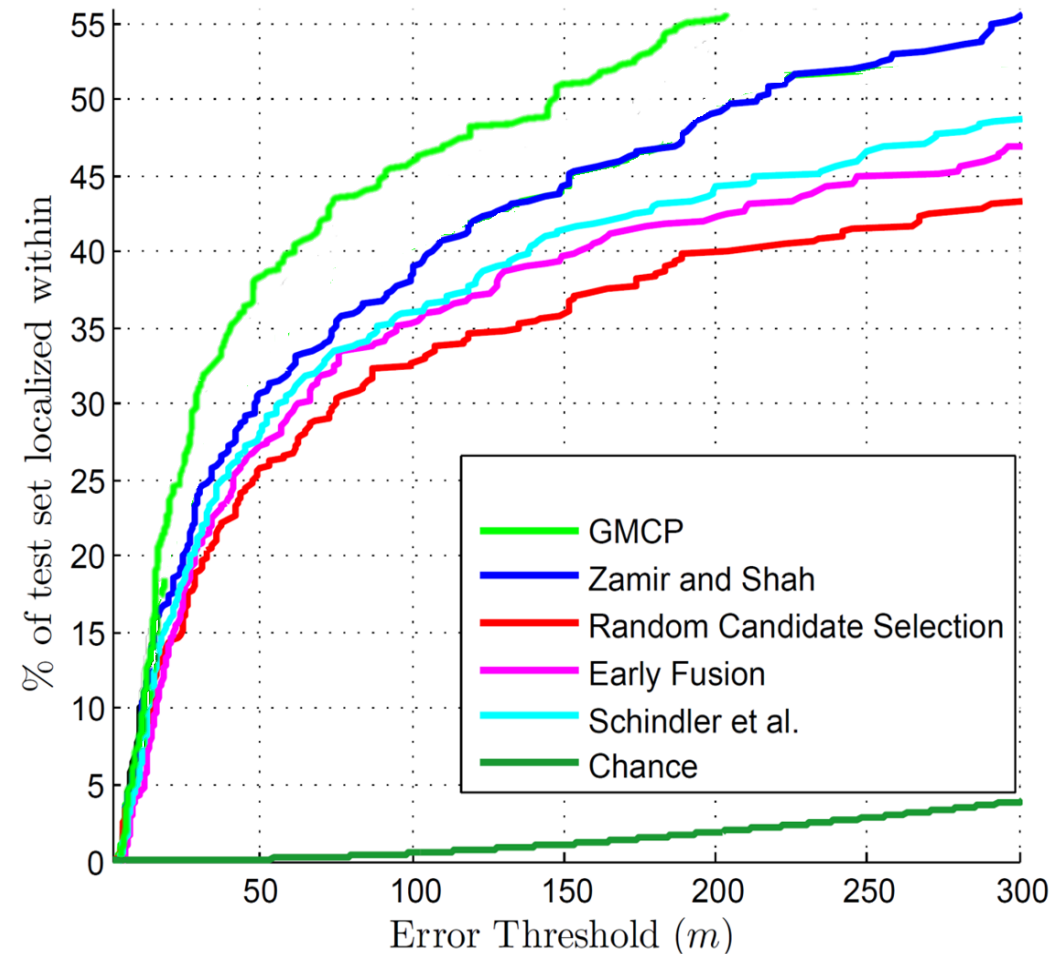


Query



Match – Error: 59.3 m

Geo-localization Results



Limitations

- GMCP selects exactly one NN per query feature; sensitive to outliers
- A very simple voting scheme
- GMCP is a binary-variable NP hard problem

Image Geo-Localization Using Constraint Dominant Sets

Eyasu Zemene, Yonatan Tariku Tesfaye, Haroon Idrees, Andrea Prati,
Marcello Pelillo, and Mubarak Shah

T-PAMI, 2017.



Eyasu Mequaanint, Qualcomm



Yonatan Tariku. , Qualcomm

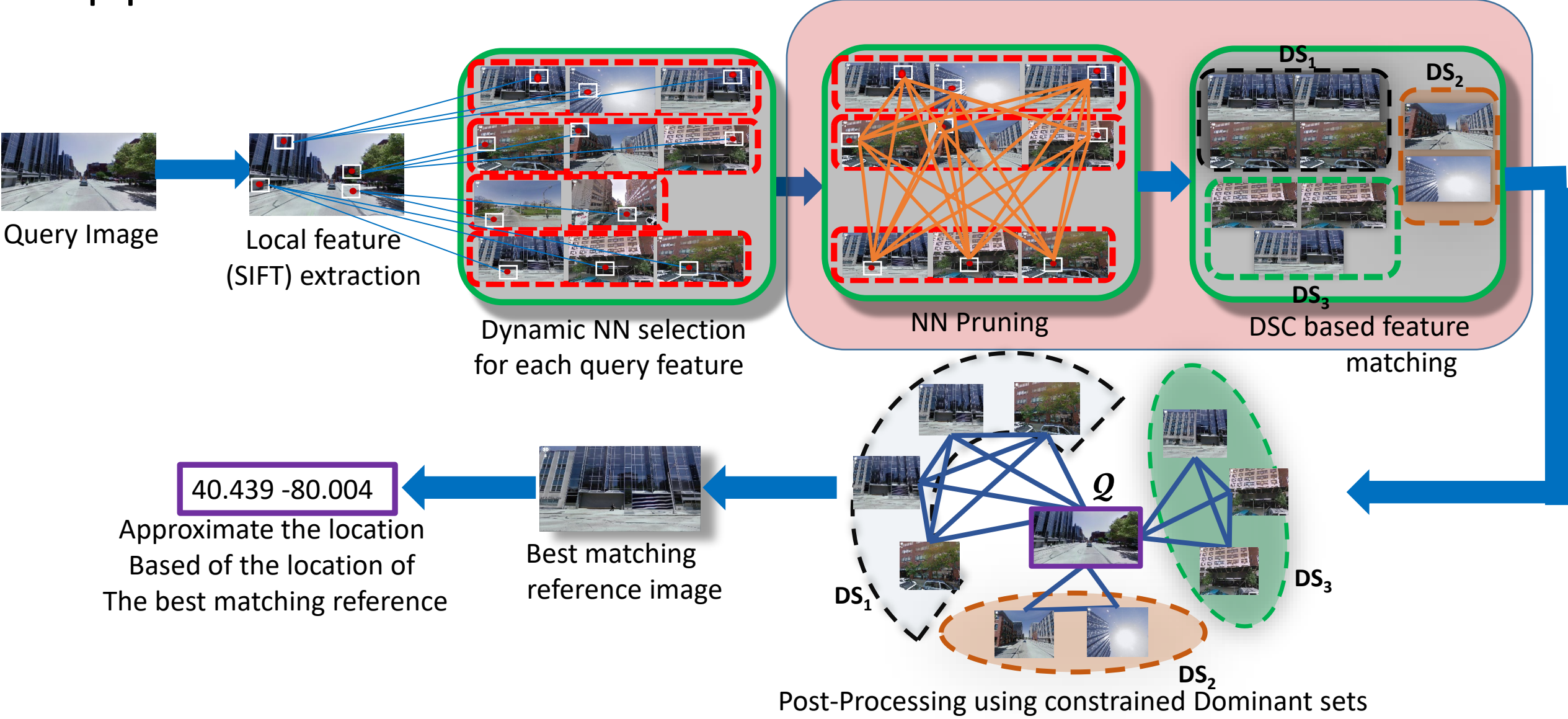


Haroon Idrees, RetailNext

Goal

- Fast
- Accurate
- Handle outliers
- Scalable to large scale

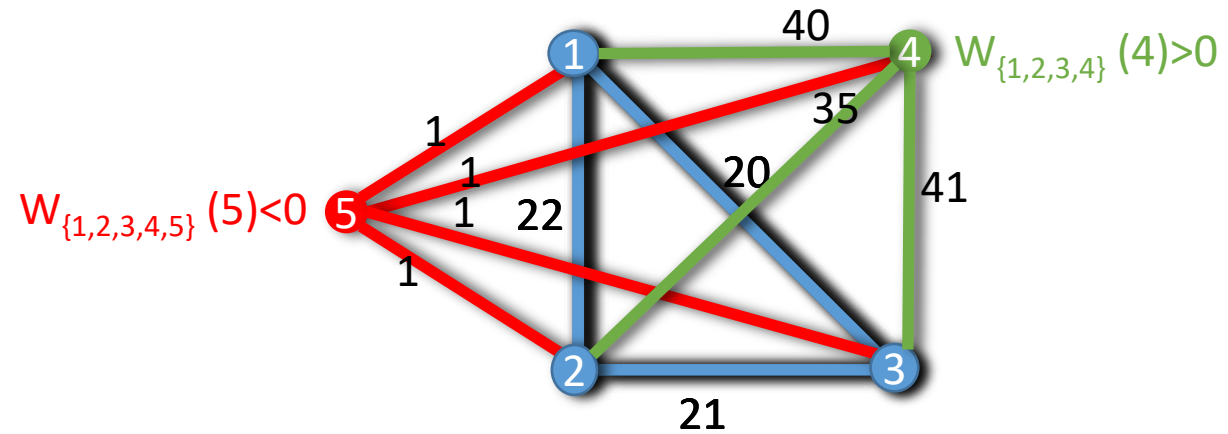
Approach



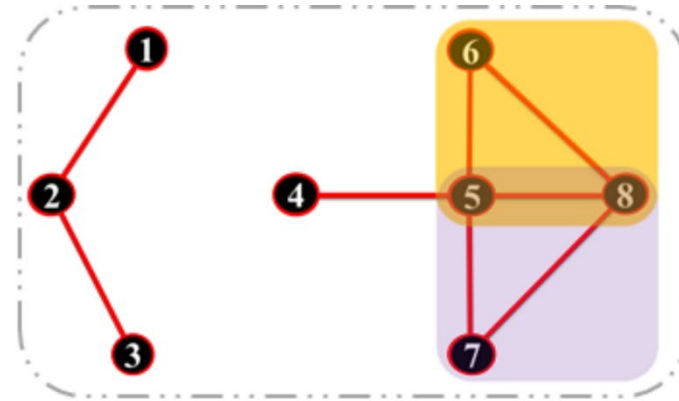
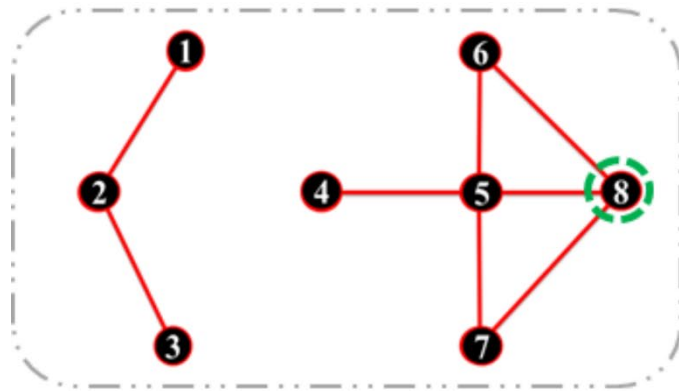
Dominant Sets

- Dominant set is an edge-weighted generalization of a clique
- Dominant set is a subset of vertices, which is
 - Coherent and
 - Compact

Dominant Sets



Constraint Dominant Sets



Dominant Sets

- Edge-weighted generalization of maximal cliques
- Given a (symmetric) affinity \mathbf{A} , consider,

$$\begin{array}{ll}\text{maximize} & f(\mathbf{x}) = \mathbf{x}' \mathbf{A} \mathbf{x} \\ \text{subject to} & \mathbf{x} \in \Delta\end{array}$$

Where $\Delta = \{\mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1 \text{ and } x_i \geq 0, \text{ for all } i = 1, \dots, n\}$

- If \mathbf{x} is a local maximizer, its support, $\sigma(\mathbf{x})$, is a dominant set
- DS's capture both *internal* and *external* coherence conditions for a cluster

M. Pavan and M. Pelillo, “Dominant sets and pairwise clustering,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 1, pp. 167–172, 2007.

Constrained Dominant Sets (CDS)

Given a query $Q \subseteq V$ and a parameter $\alpha > 0$, define the following parameterized family of quadratic program:

$$\begin{aligned} & \text{maximize } f_Q^\alpha = x^T (A - \alpha I_Q) x \\ & \text{Subject to } x \in \Delta \end{aligned} \tag{1}$$

Where I_Q is the diagonal matrix whose **diagonal elements are set to 1** in correspondence to the vertices contained in $V \setminus Q$ and to 0 otherwise.

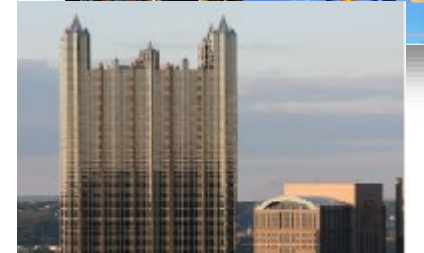
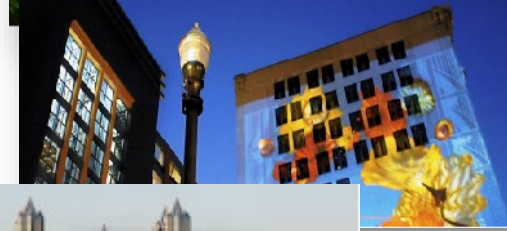
Experimental Results

Dataset I: Orlando & Pittsburgh:

- Reference images:
 - 102K Google street view images from **Pittsburgh, PA** and **Orlando, FL**
- Test Set:
 - 644 GPS-Tagged unconstrained images
 - Downloaded From Flickr, Panoramio, Picasa, ...

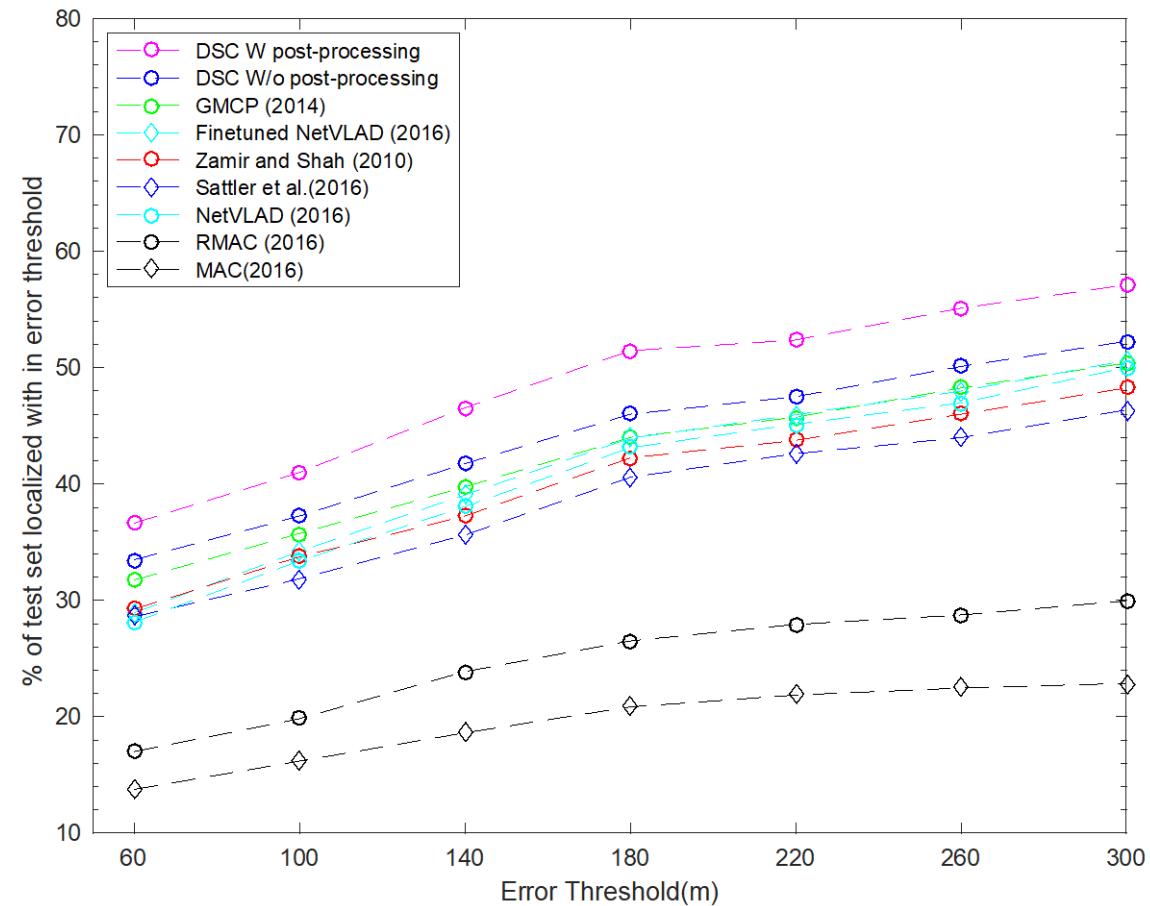
Dataset II: *WorldCities* Dataset: (NEW)*

- **Reference images** (300K Google street view images):
- 14 different cities from different parts of the world:
 - **USA:** Los Angeles, Phoenix, Houston, San Diego, Las Vegas, Dallas, Chicago
 - **Australia:** Sydney and Melbourne
 - **Europe:** Amsterdam, Frankfurt, Rome, Milan and Paris
- **Test Set**
 - 500 GPS-Tagged unconstrained images
 - Downloaded From Flickr, Panoramio, Picasa...



Overall Results

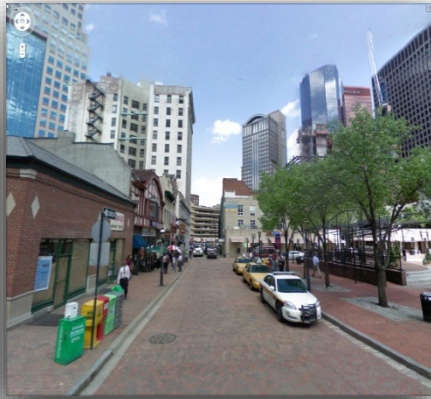
- Dataset 2: WorldCities (14 different cities from Europa, North America, Australia)



Qualitative Image Localization Results



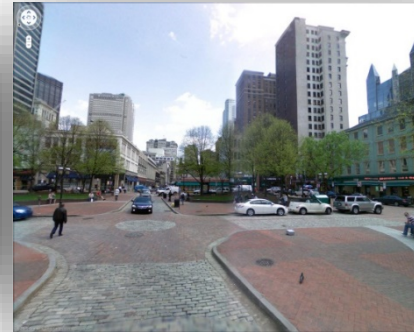
Query



Match – Error: 70.01 m



Query



Match – Error: 5.4 m



Query



Match – Error: 10.4 m



Query



Match – Error: 7.5 m

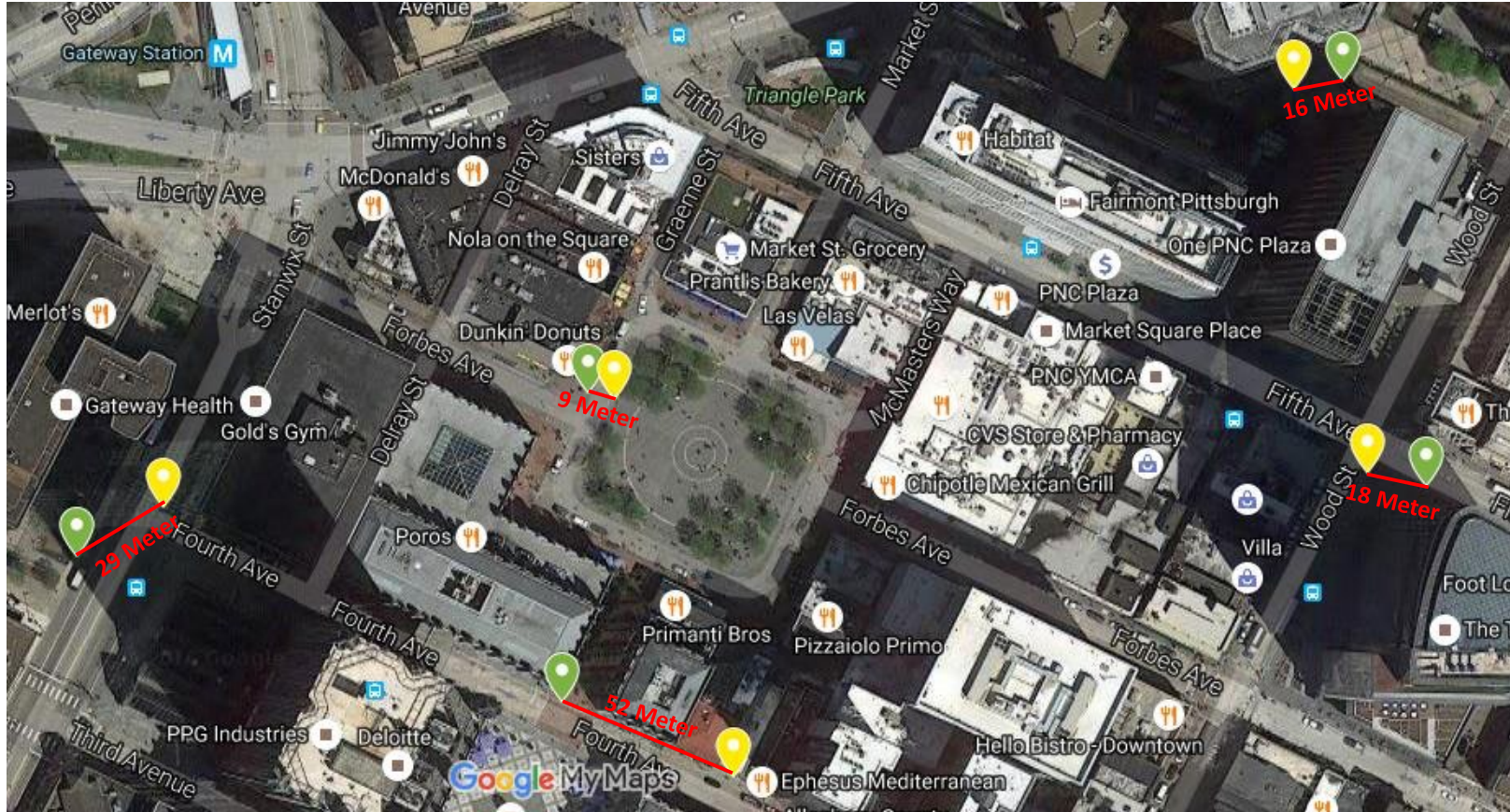


Query



Match – Error: 62.7 m

Qualitative Image Localization Results

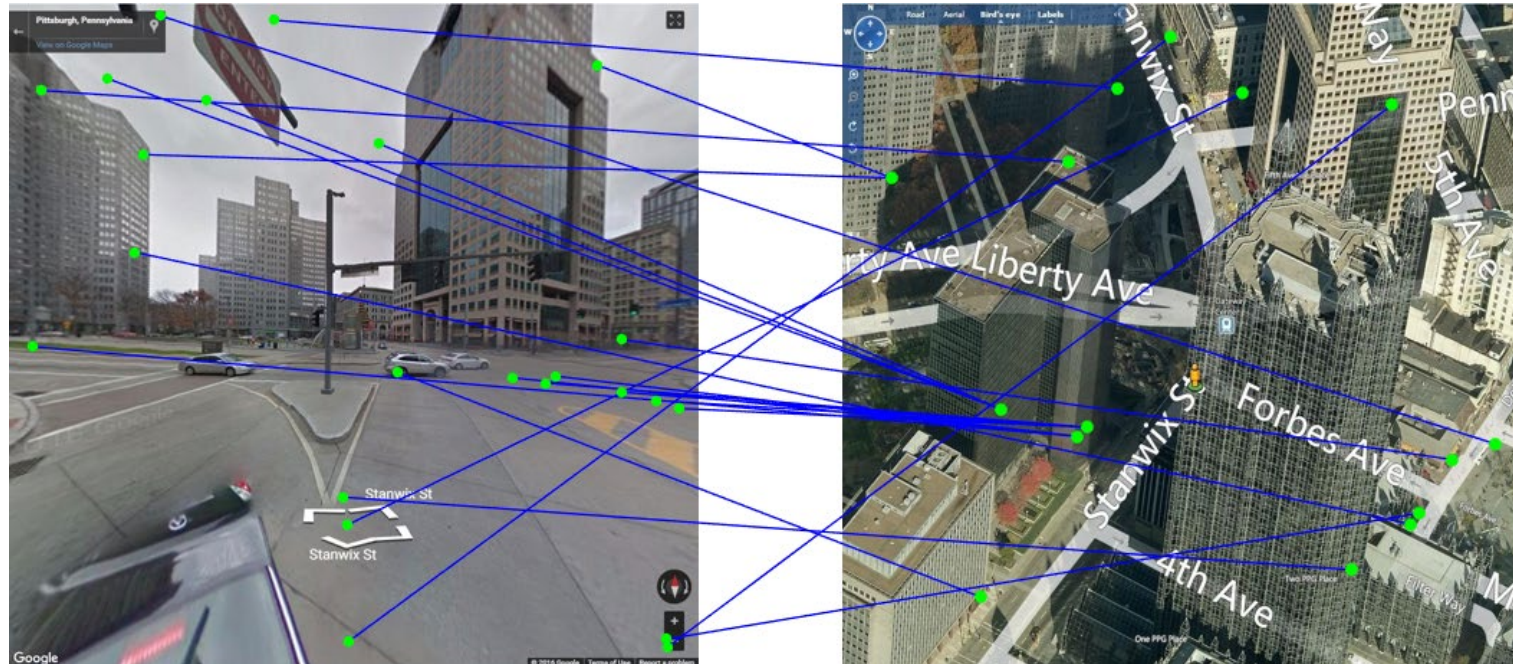


Contents

- Pixel-Wise Geo-localization
 - Geodetic Alignment of Aerial Video Frames
- Image-Based Geo-Localization
 - Same View (Street-View to Street-View)
 - Generalized Maximum Clique (PAMI, 2014)
 - Constraint Dominant Sets (PAMI, 2017)
 - **Cross-View Geo-Localization**
 - Bird's Eye-View to Street View (CVPR, 2017)
 - Aerial to Ground View (ICCV, 2019)

Cross-View Challenges

- Images from different viewpoints are visually different
- The images may be captured with different lighting
- The mapping from one viewpoint to the other may be complex
- Traditional features e.g. SIFT, HOG etc. may be very different



Retrieval Features

- Local Features (SIFT) (same view)
- **Buildings Features** (cross view)

Cross-View Image Matching for Geo-localization in Urban Environments

Yicong Tian, Chen Chen, Mubarak Shah
Center for Research in Computer Vision,
University of Central Florida
CVPR-2017

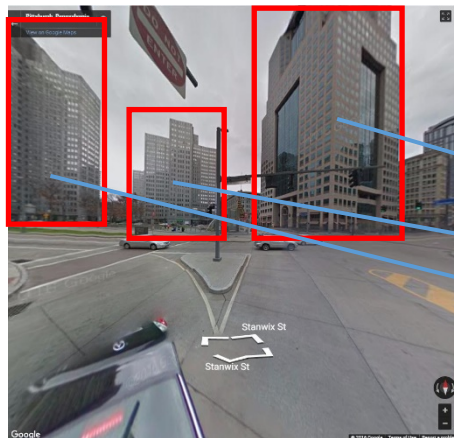
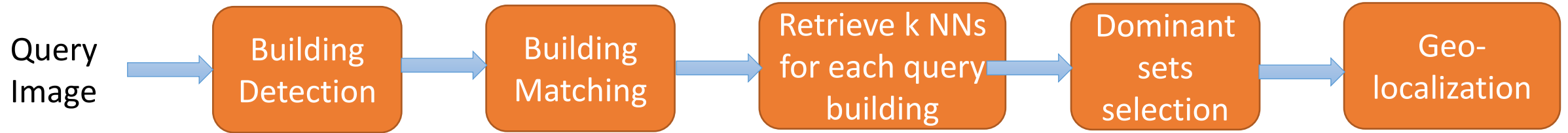


Yicong Tian, Google

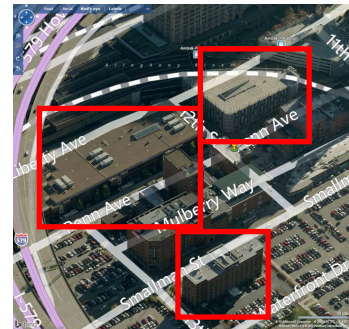


Chen Chen, UNC Charlotte

Proposed Approach



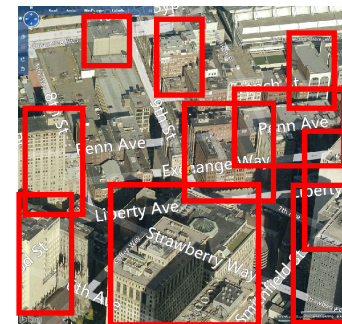
Query image
(street view)



...



...



Reference
images
(bird's eye
view)

Data Collection

- A new dataset of street view and bird's eye view image pairs
 - Pittsburg
 - Orlando
- Generated a list of GPS coordinates along streets
- For each GPS location
 - Four pairs with different headings
 - Utilized [DualMaps](#)
 - Automatically saved screenshots
- Annotations
 - By four undergraduates and high school students
 - ~300 hours work in total

Data Collection - Pittsburgh

- 1,586 GPS locations



Data Collection - Pittsburg

- Example image pair with annotations



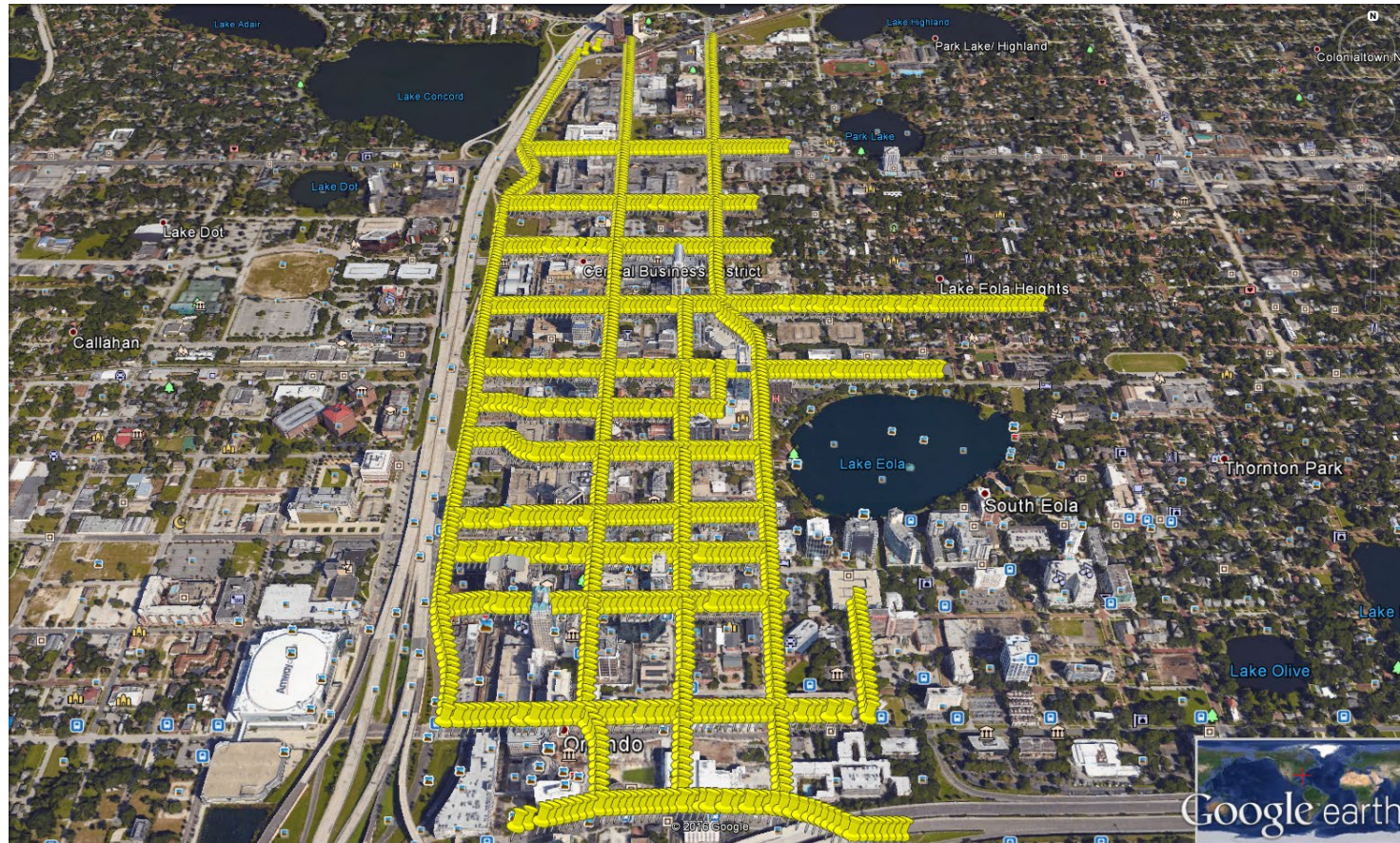
Data Collection - Pittsburg

- Example image pair with annotations



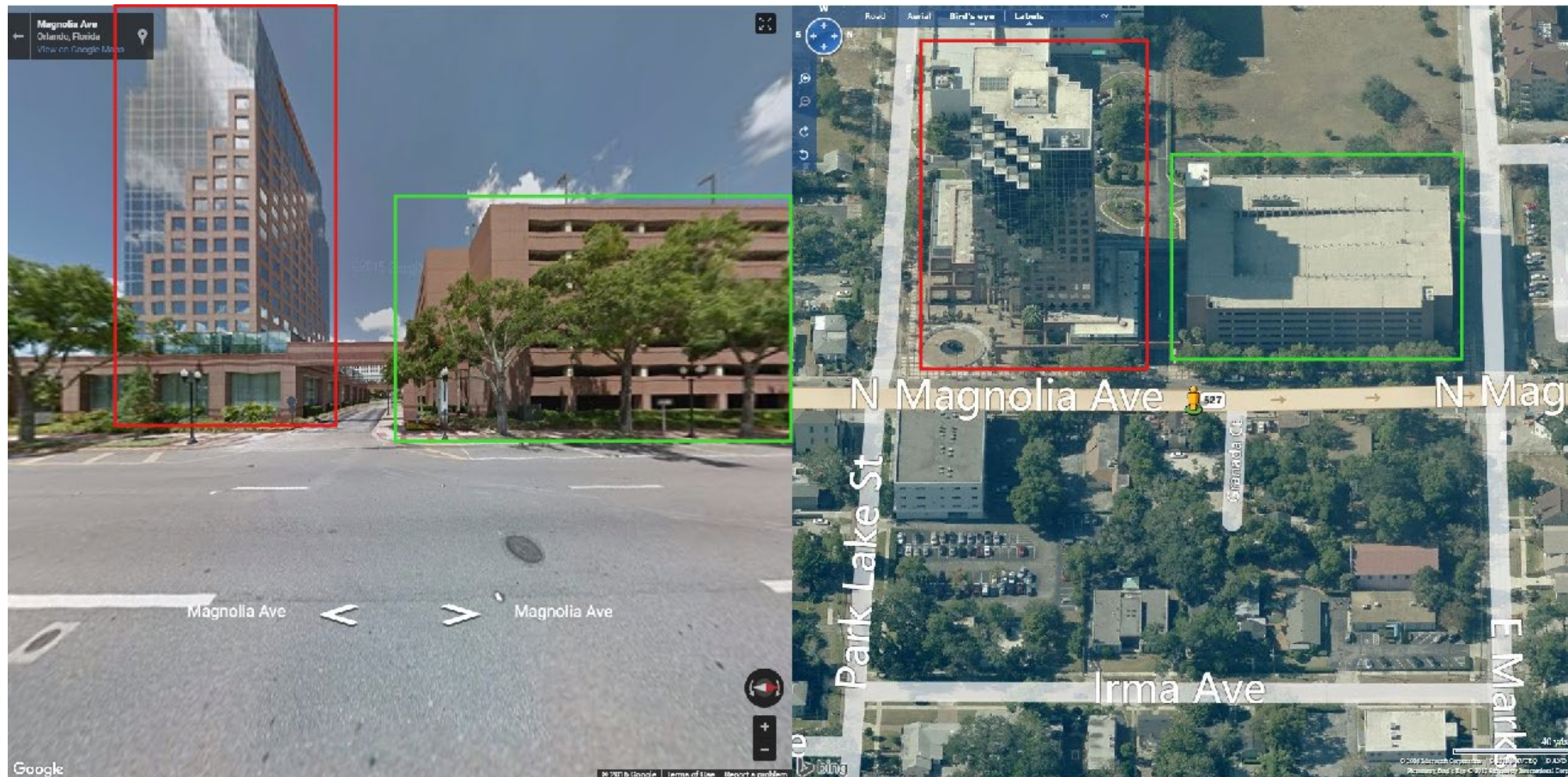
Data Collection - Orlando

- 1,324 GPS locations



Data Collection - Orlando

- Example image pair with annotations

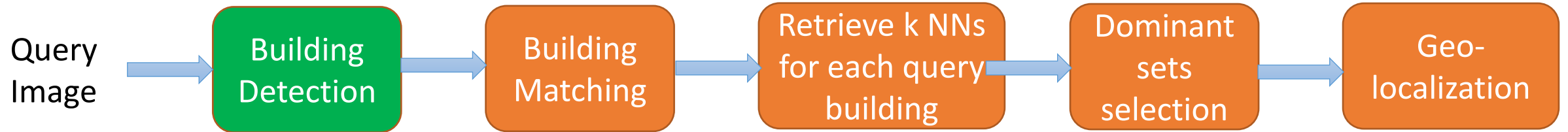


Data Collection - Orlando

- Example image pair with annotations



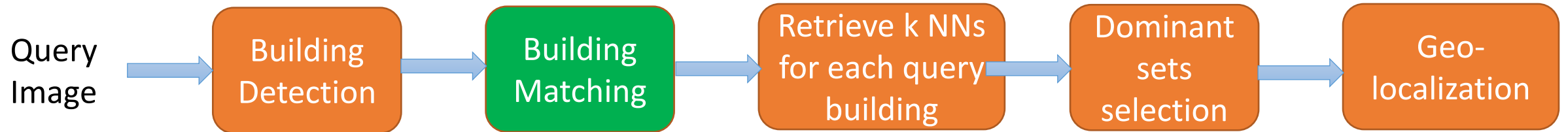
Proposed Approach



Building Detection (Faster R-CNN)

- Street view images
 - Training: 6,682 images
 - With 15,648 annotated boxes
 - Test: 1,903 images
- Bird's eye view images
 - Training: 6,968 images
 - With 39,511 annotated boxes
 - Test: 1,916 images
- Each model takes 10 hours to train

Proposed Approach



Building Matching

- Train a Siamese network
- In the learned feature space,
 - Matching image pairs are close to each other
 - Unmatched image pairs are far apart
- Contrastive loss

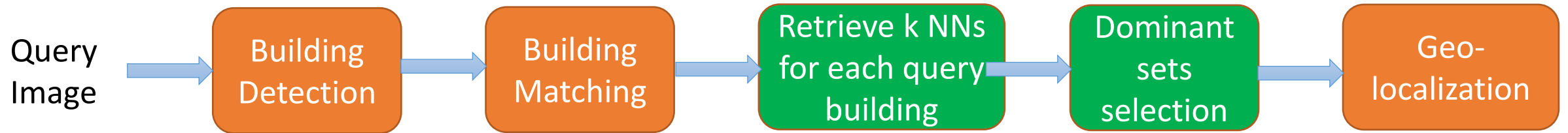
$$\mathcal{L}(x, y, l) = \frac{1}{2}lD^2 + \frac{1}{2}(1 - l) \max(0, (m - D^2))$$

l: label (1 or 0)

D^2 : square of Euclidean distance between features

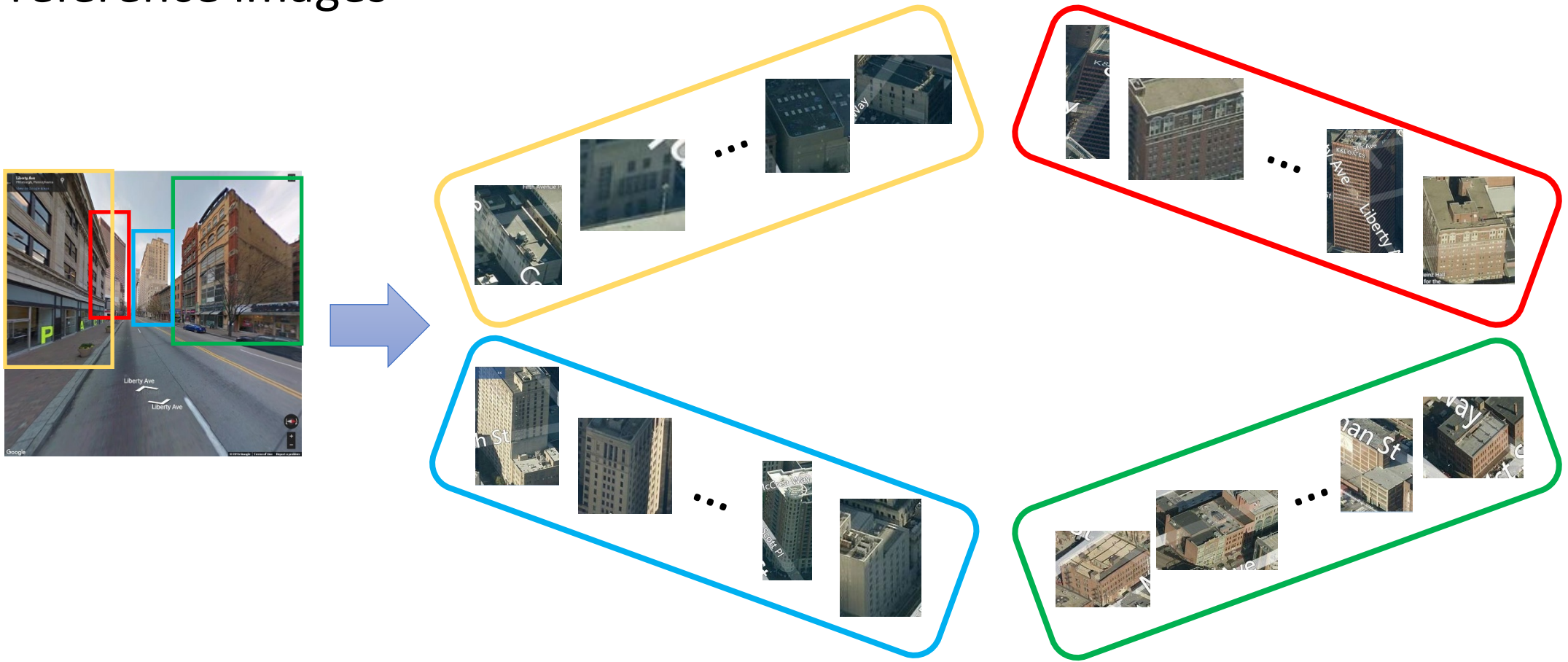
m: margin

Proposed Approach



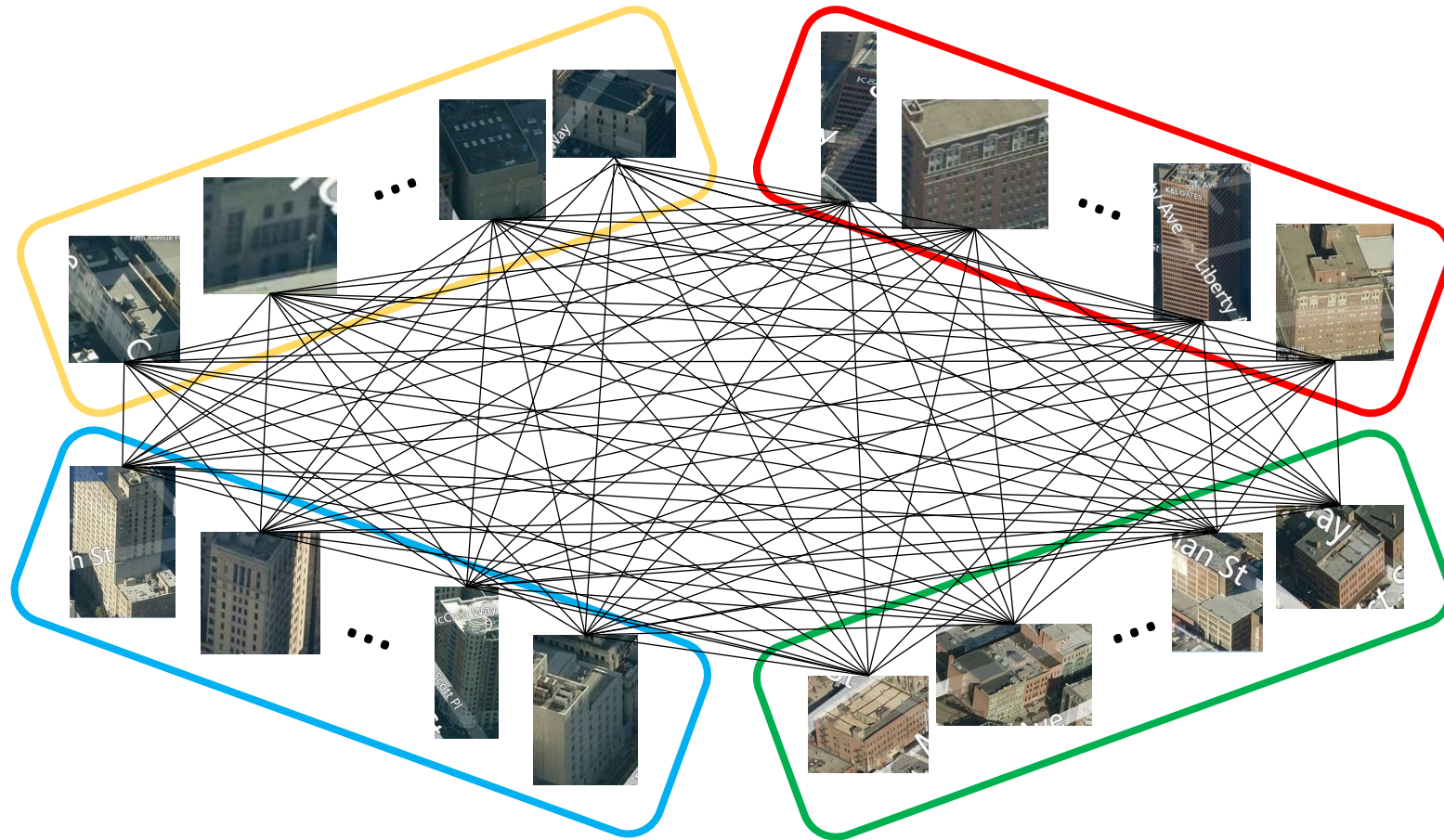
Geo-localization Using Dominant Sets

- For each building in the query image, select k nearest neighbors from reference images



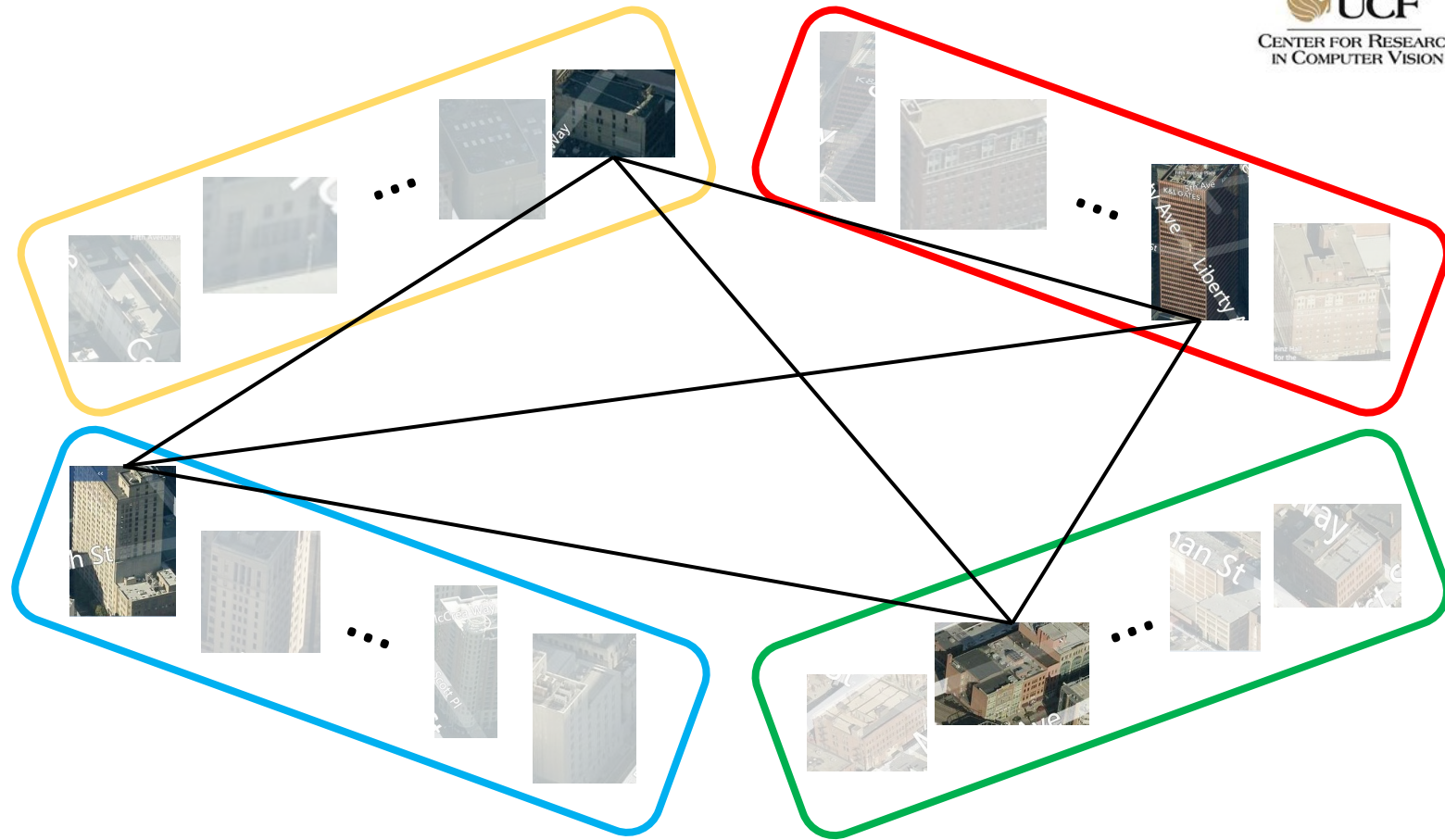
Geo-localization Using Dominant Sets

- Build a graph $G = (V, E, \omega)$ using selected reference buildings



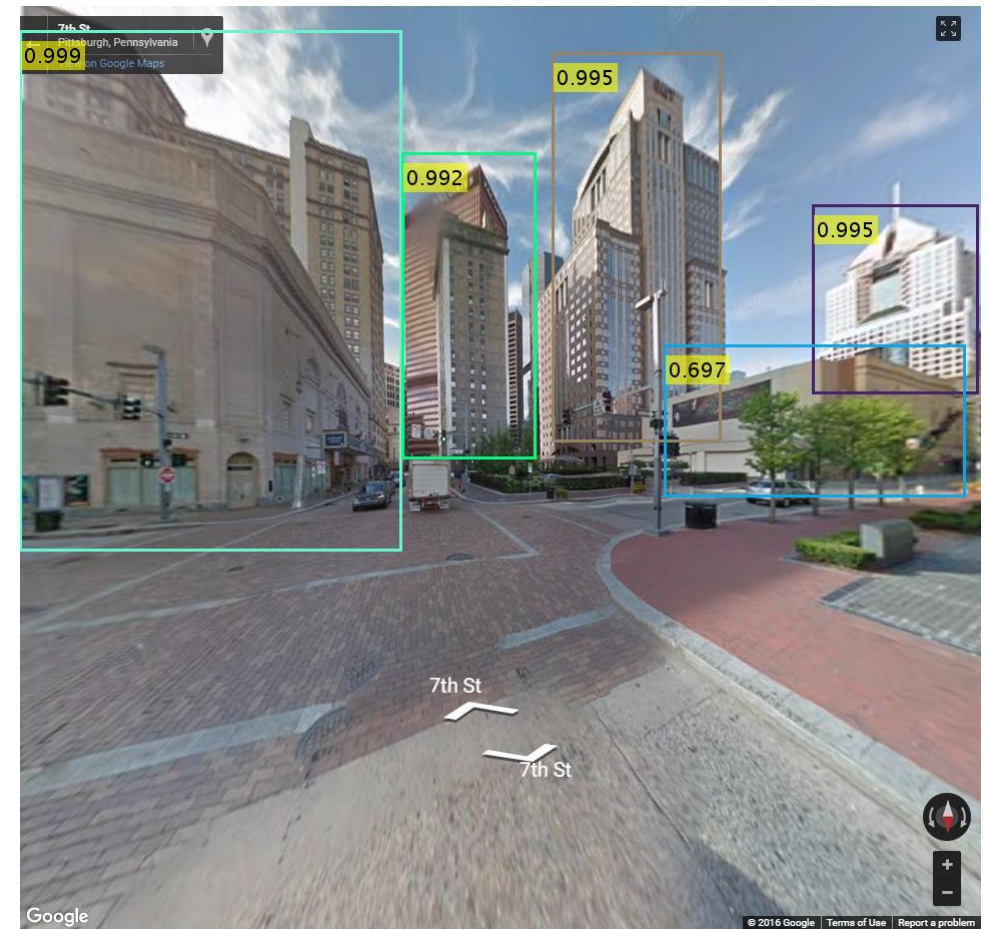
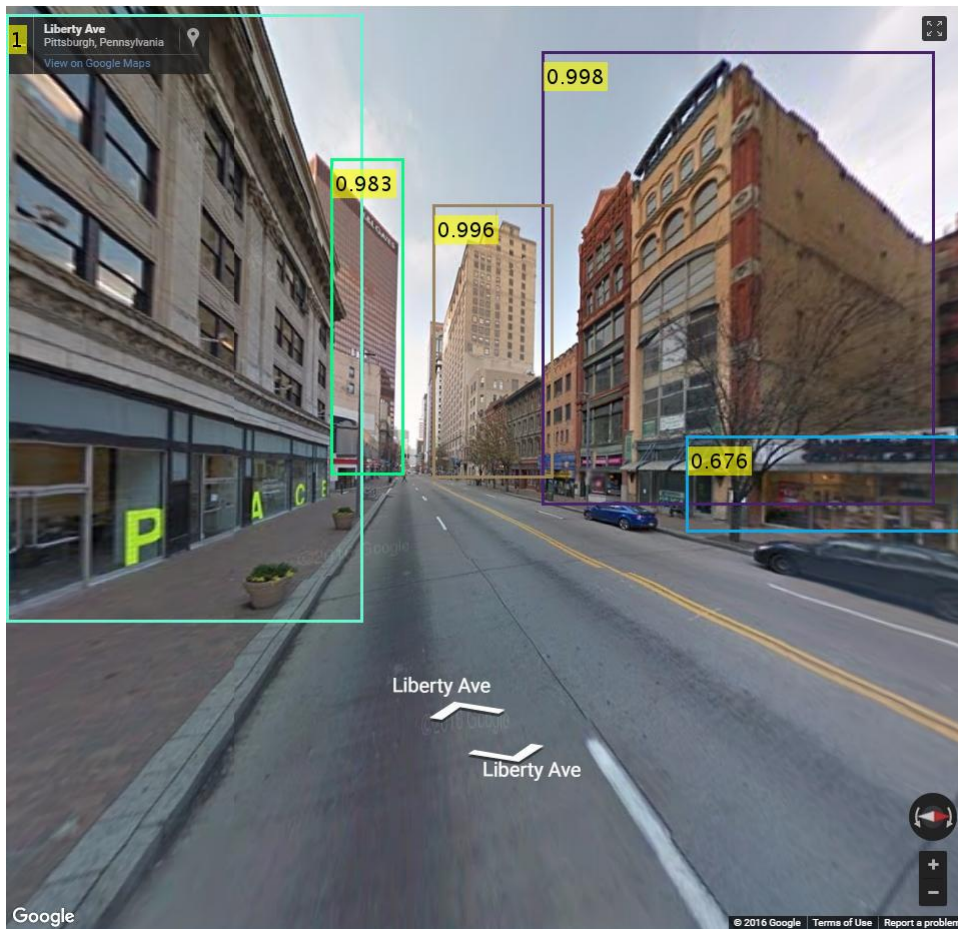
Geolocalization

- The nodes in dominant set form a coherent set
- At most one node is selected from each cluster

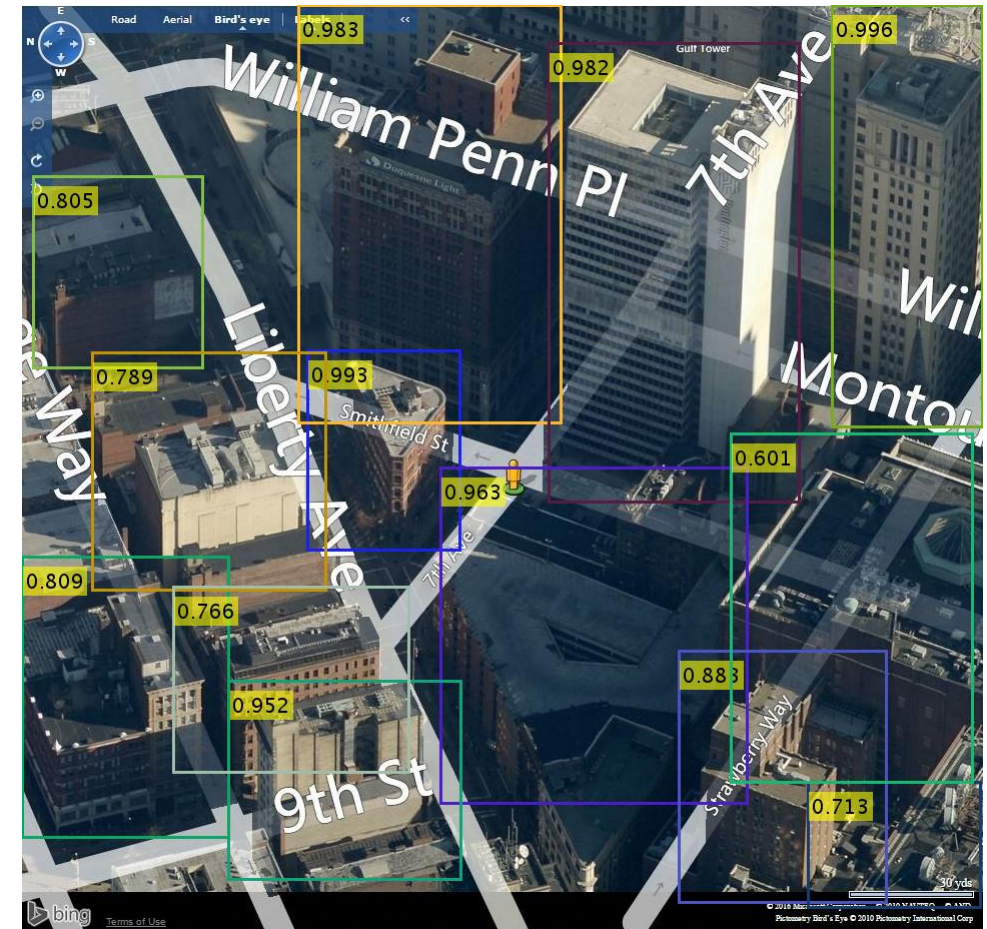
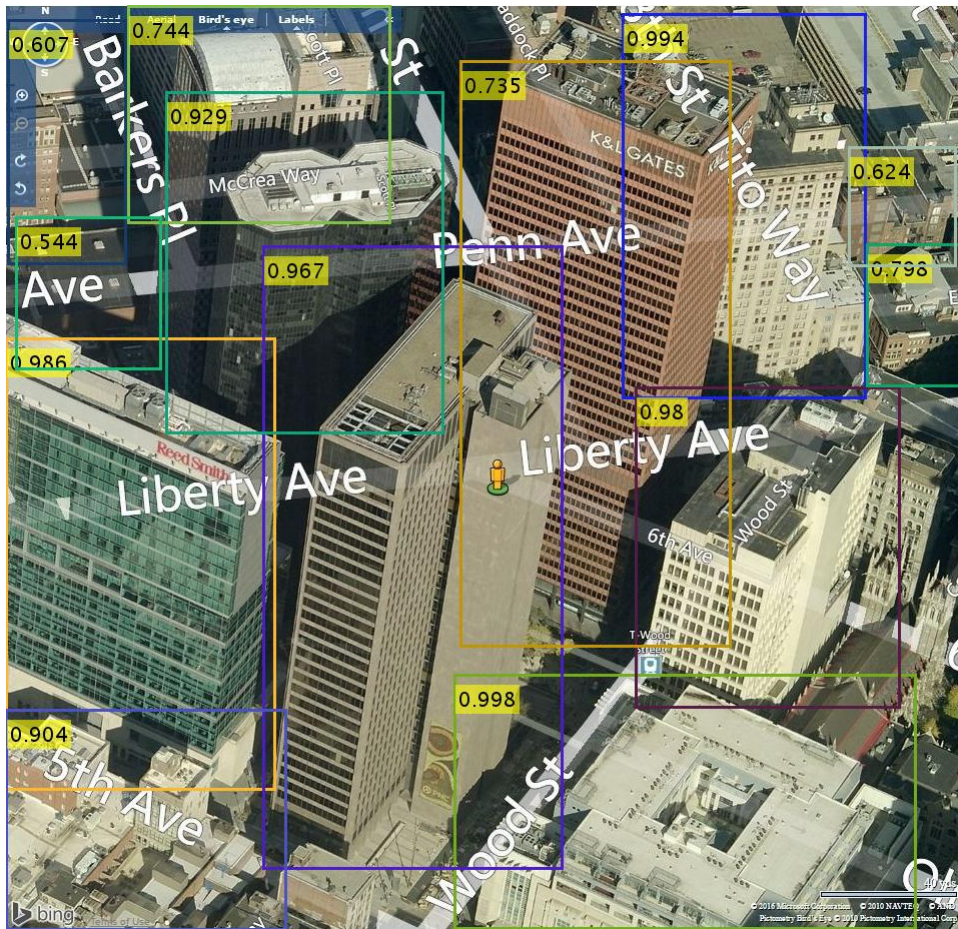


Experimental Results

Building detection in street view images

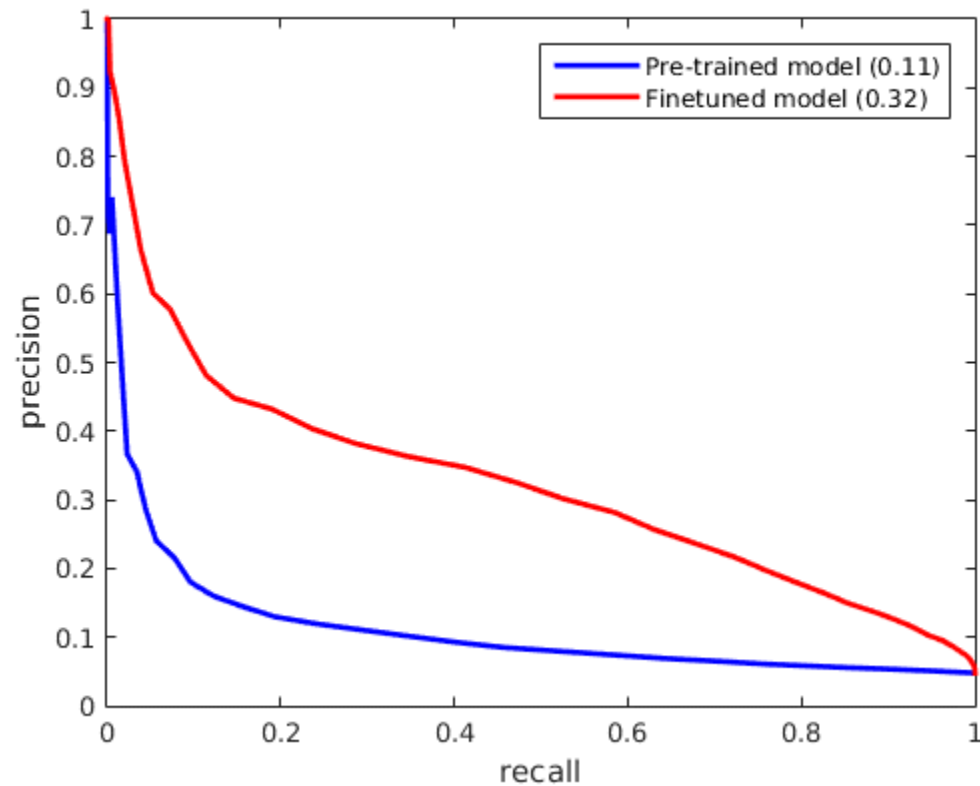


Building detection in bird's eye view images



Building matching

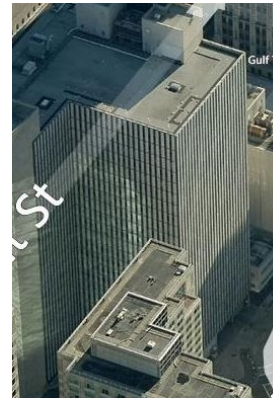
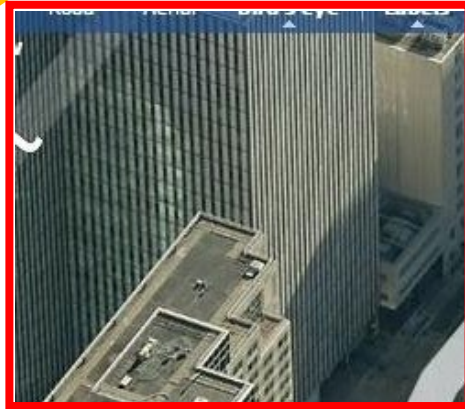
- Precision-recall curve on test image pairs



Building matching



Query image
(street view)

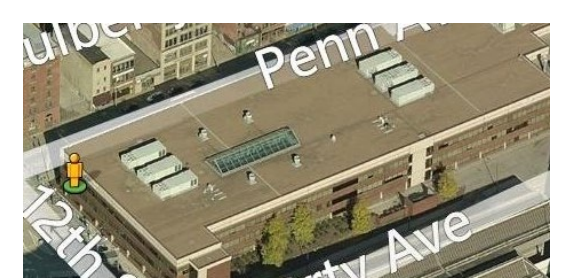
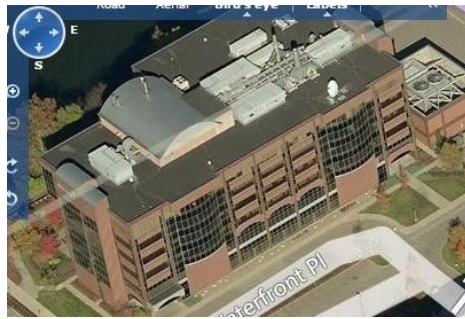


From ~40,000 candidate image patches (Buildings)
(bird's eye view)

Building matching



Query image
(street view)

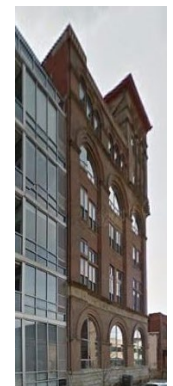
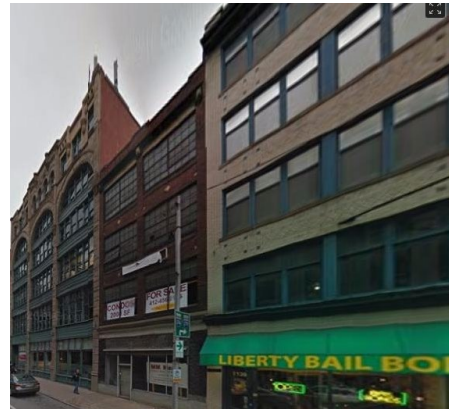
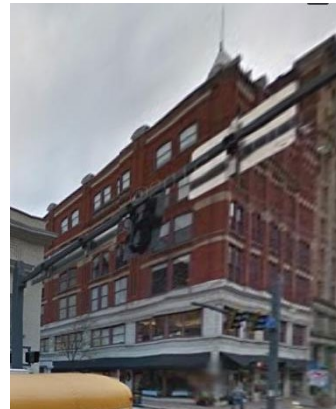
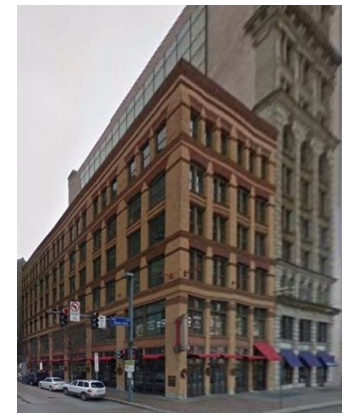
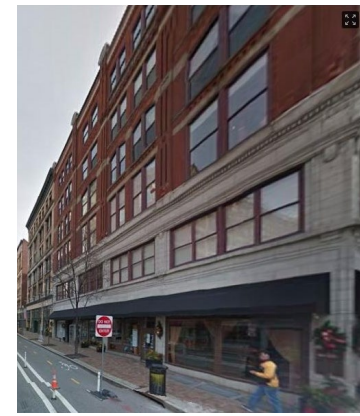
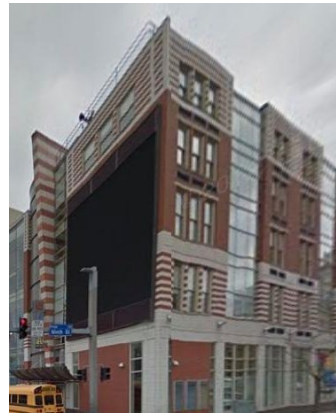
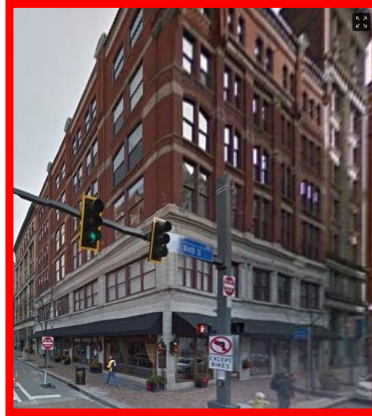


From ~40,000 candidate image patches (Buildings)
(bird's eye view)

Building matching



Query image
(bird's eye view)



From ~10,000 candidate image patches (Buildings)
(street view)

Building matching



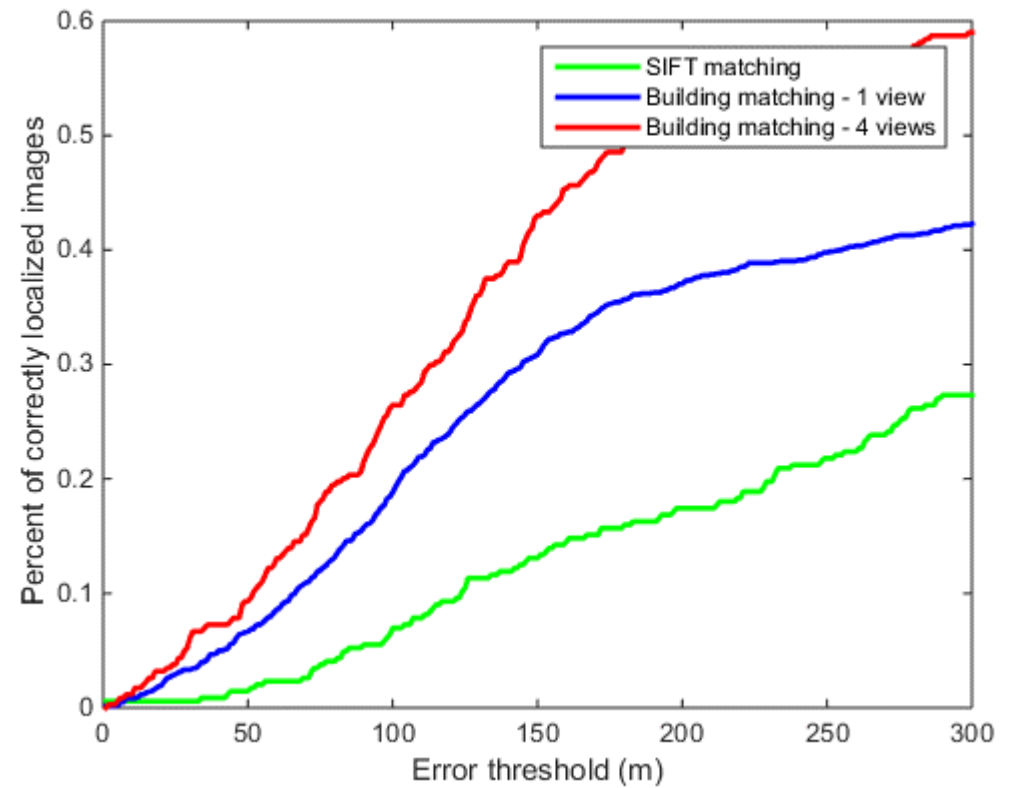
**Query image
(bird's eye view)**



**From ~10,000 candidate image patches (Buildings)
(street view)**

Geo-localization

- Query image: street view
- Reference images: bird's eye view
- k: 100



Summary

- Geo-localization Using Cross View Image Matching
 - Detect Buildings
 - Match Buildings
 - Retrieve k-nearest neighbors for each query
 - Dominant Set Selection

Cross-View Image Matching for Geo-localization in Urban Environments

Yicong Tian, Chen Chen, Mubarak Shah

Center for Research in Computer Vision, University of Central Florida

CVPR-2017

Retrieval Features

- Local Features (SIFT)
- Building Features
- **Global Image Features**

Bridging the Domain Gap for Ground-to-Aerial Image Matching



Krishna Regmi

Krishna Regmi & Mubarak Shah
University of Central Florida

https://www.youtube.com/watch?v=gmAhQXCCEQ&list=PLd3hISJsX_IkSnnrMtzsMHI1q6vimipvp&index=1

Introduction

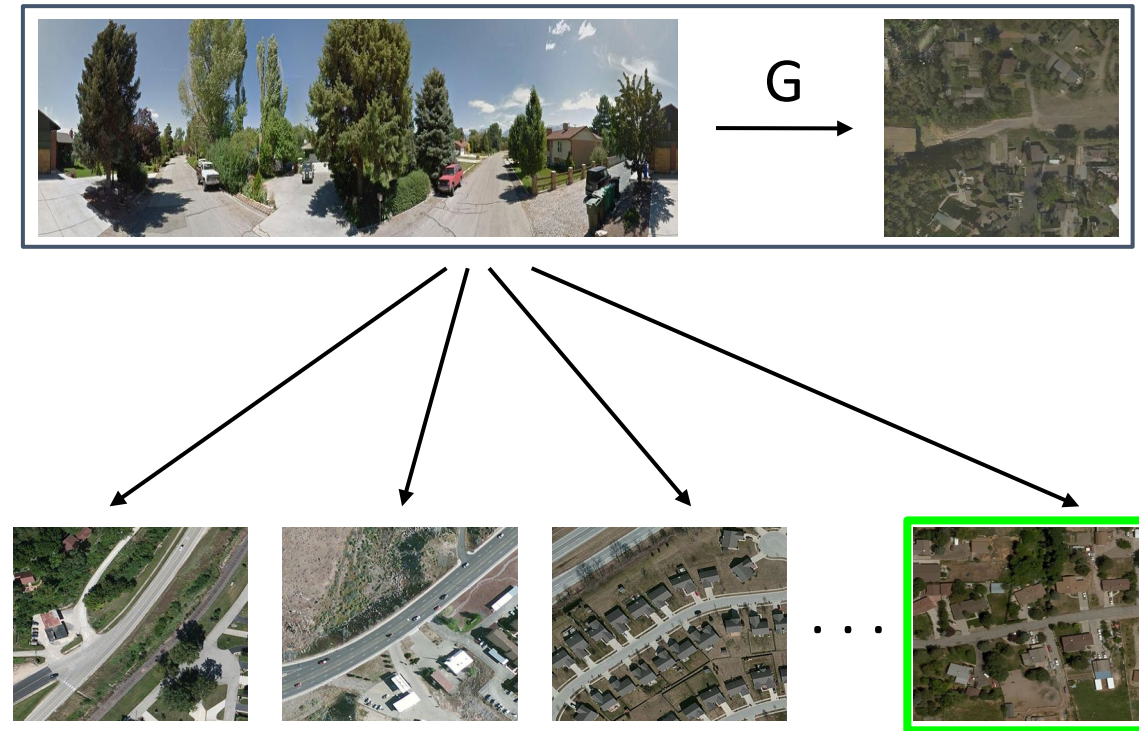
Cross-view image matching

Drastic viewpoint differences

Bridge the domain gap - use GANs
(synthesize target-view images)

Joint Feature Learning and Feature fusion

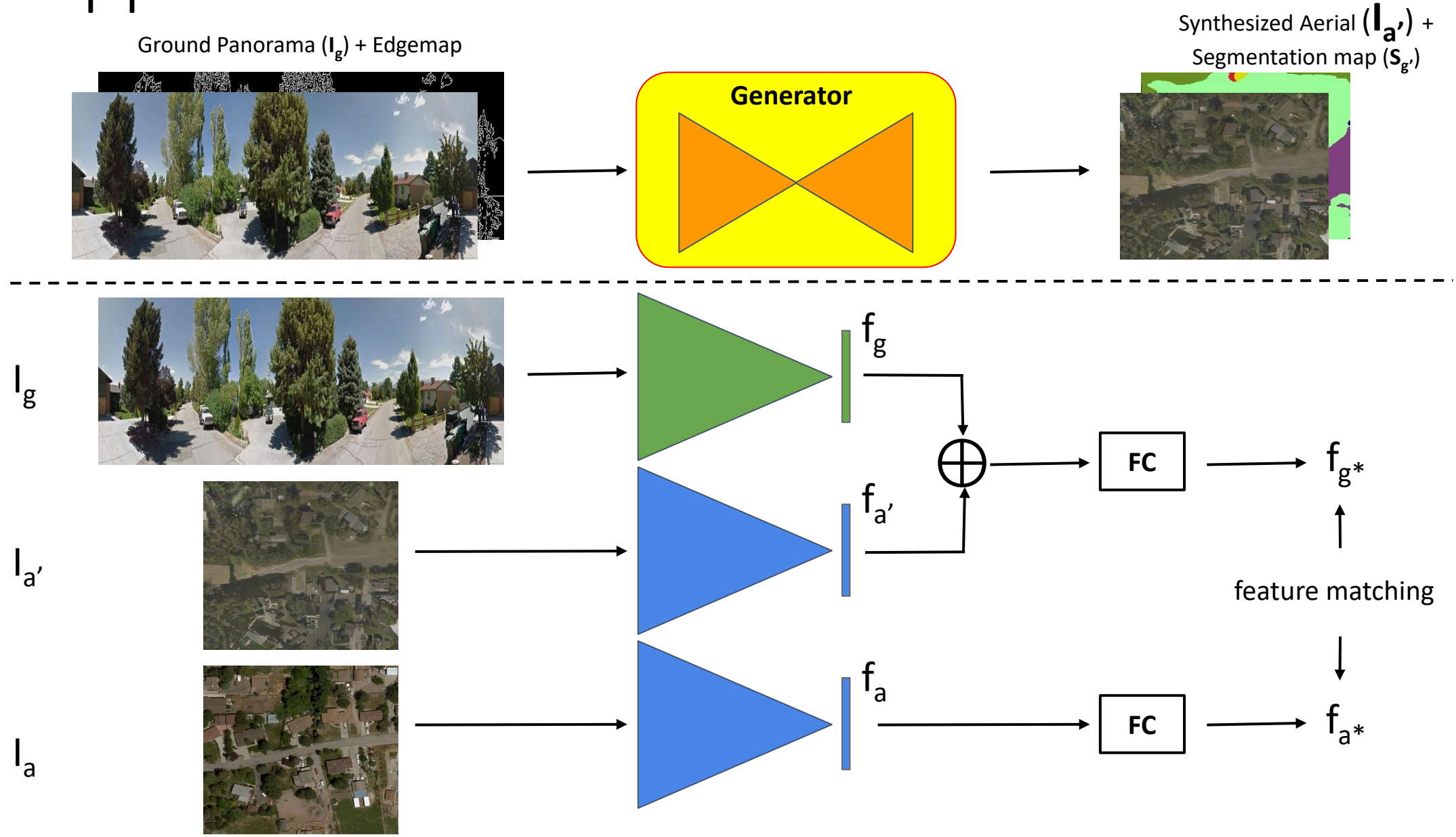
Multi-scale feature aggregation



Our Approach

Cross-view image synthesis followed by
Joint Feature Learning and Feature Fusion

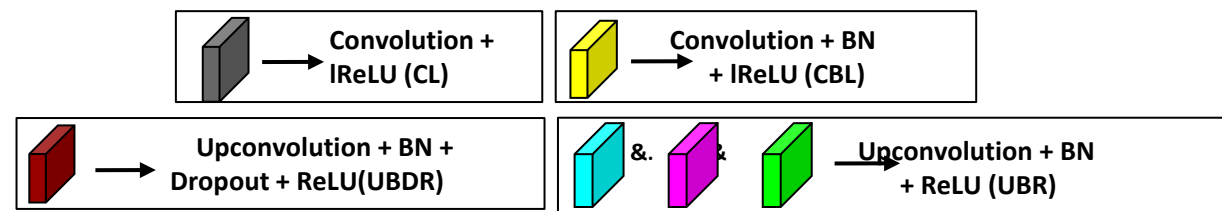
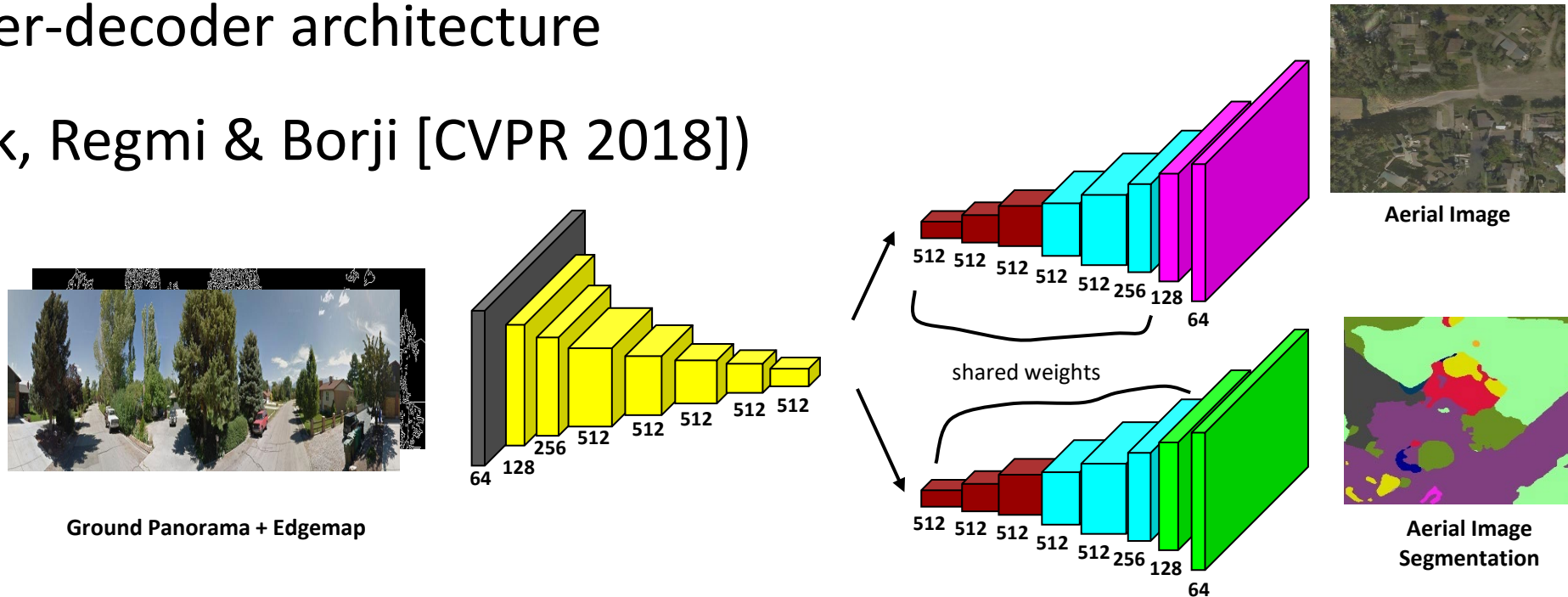
Our Approach



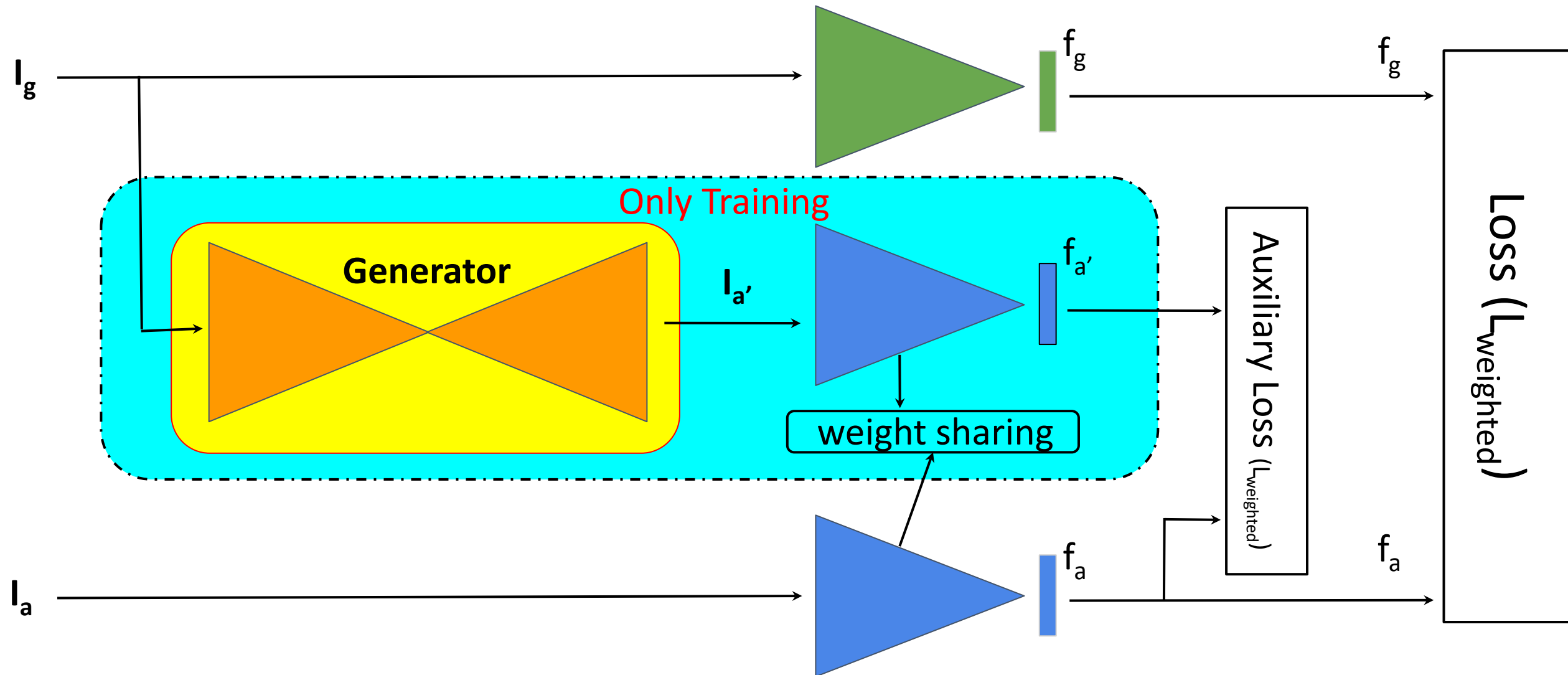
Generator Architecture

Encoder-decoder architecture

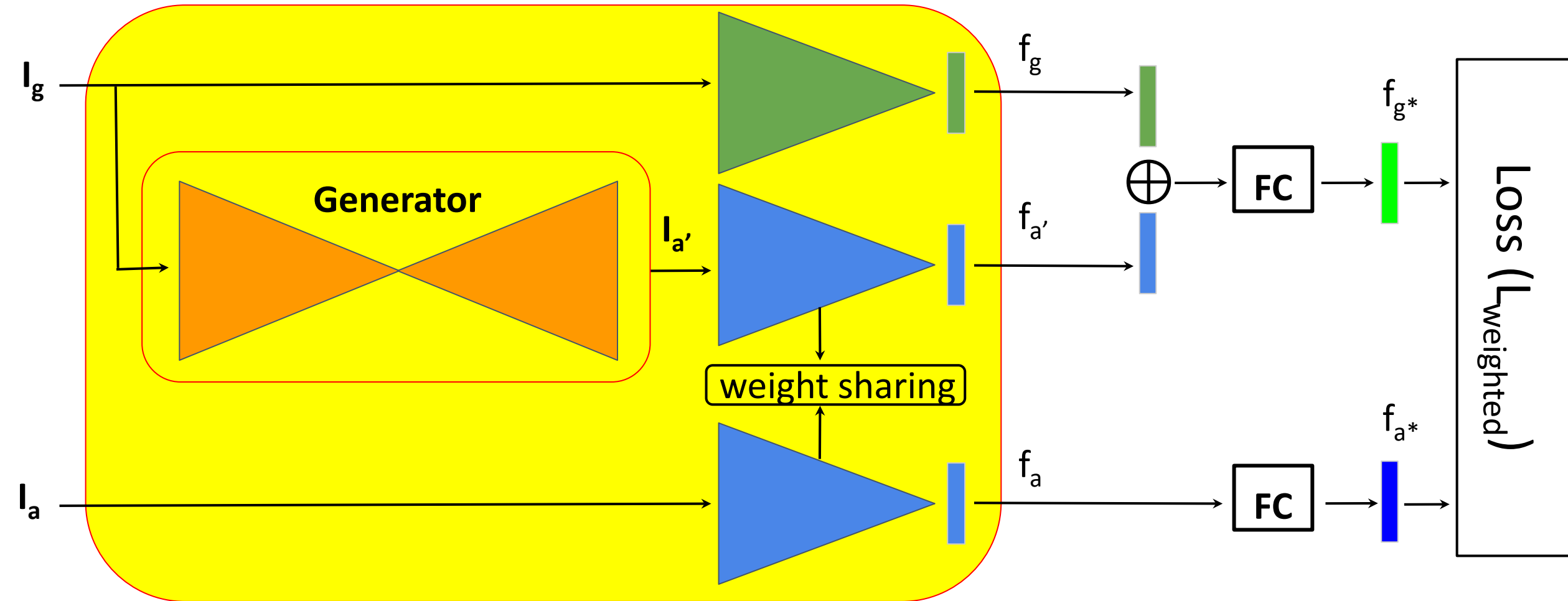
(X-Fork, Regmi & Borji [CVPR 2018])



Joint Feature Learning



Feature Fusion



Loss Functions

Triplet Loss:
$$L_{triplet} = \max(0, m + d_p - d_n)$$
$$= \max(0, m + d) \quad \text{where, } d = d_p - d_n$$

Soft-Margin Triplet Loss:
$$L_{soft} = \ln(1 + e^d)$$

Weighted Soft-Margin Triplet Loss:
$$L_{weighted} = \ln(1 + e^{\alpha d})$$

Loss for Joint Feature Learning:

$$L_{joint} = \lambda_1 L_{weighted}(I_g, I_a) + \lambda_2 L_{weighted}(I_{a'}, I_a)$$

Datasets

Satellite - ground panorama pairs

CVUSA Dataset:

Train/Test: 35,532/8,884 pairs

Covers rural areas

UCF OP Dataset: (Orlando Pittsburgh)

Train/Test: 1,910/722 pairs

Newly collected

Covers urban areas

GPS info available

Results: CVUSA Dataset

Quantitative Results (CVUSA Dataset)

Method	Top-1	Top-10	Top-1%
Two-stream baseline ($I_{a'}$, I_a)	10.23%	35.10%	72.58%
Two-stream baseline (I_g , I_a)	18.45%	48.98%	82.94%
Joint Feat. Learning ($I_{a'}$, I_a)	14.31%	48.75%	86.47%
Joint Feat. Learning (I_g , I_a)	29.75%	66.34%	92.09%
Feature Fusion	48.75%	81.27%	95.98%
Workman et al. [41]	-	-	34.3%
Zhai et al. [46]	-	-	43.2%
Vo and Hays [39]	-	-	63.7%
CVM-Net-I [18]	22.53%	63.28%	91.4%
CVM-Net-II [18]	11.18%	43.51%	87.2%

Recall Accuracy (CVUSA Dataset)

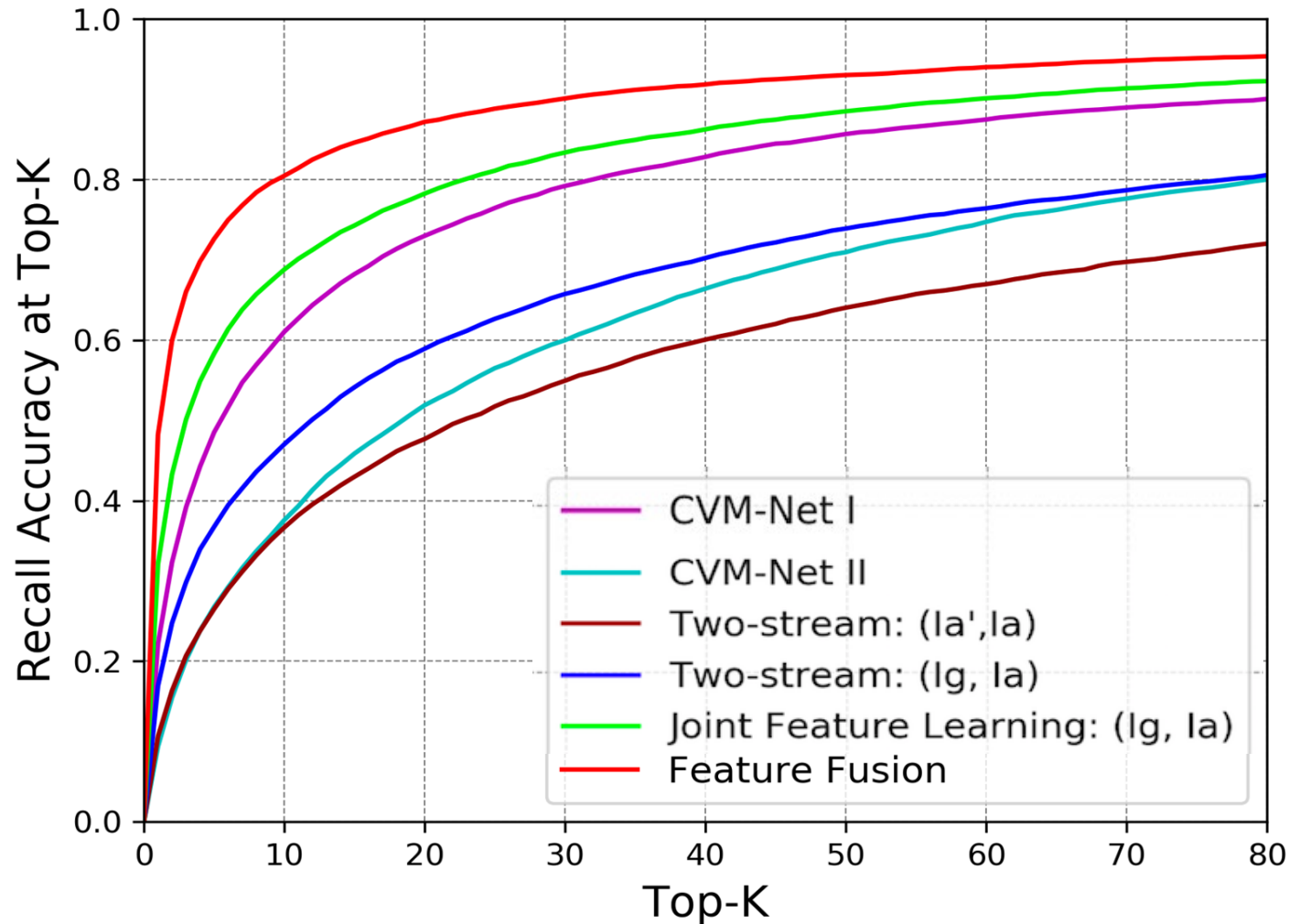


Image Retrieval (CVUSA Dataset)

Ground Query

Synthesized
Aerial

Top matches (top 1 – top 5 from left to right)



Failure Cases

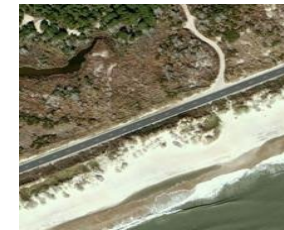
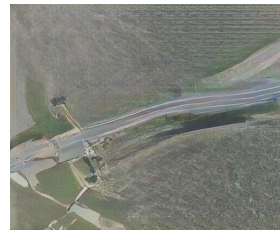
Ground Query

Synthesized Aerial

Top-1 match

Correct match
(Ground Truth)

Retrieved
position



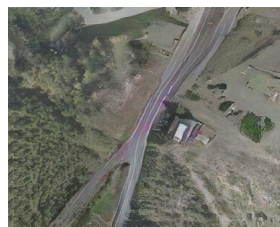
13



37



532



1700

Results: OP Dataset

Retrieval Performance (OP Dataset)

Two-stream (I_g, I_a)	Joint Feat. Learning	Feature Fusion
30.61%	38.36%	45.57%

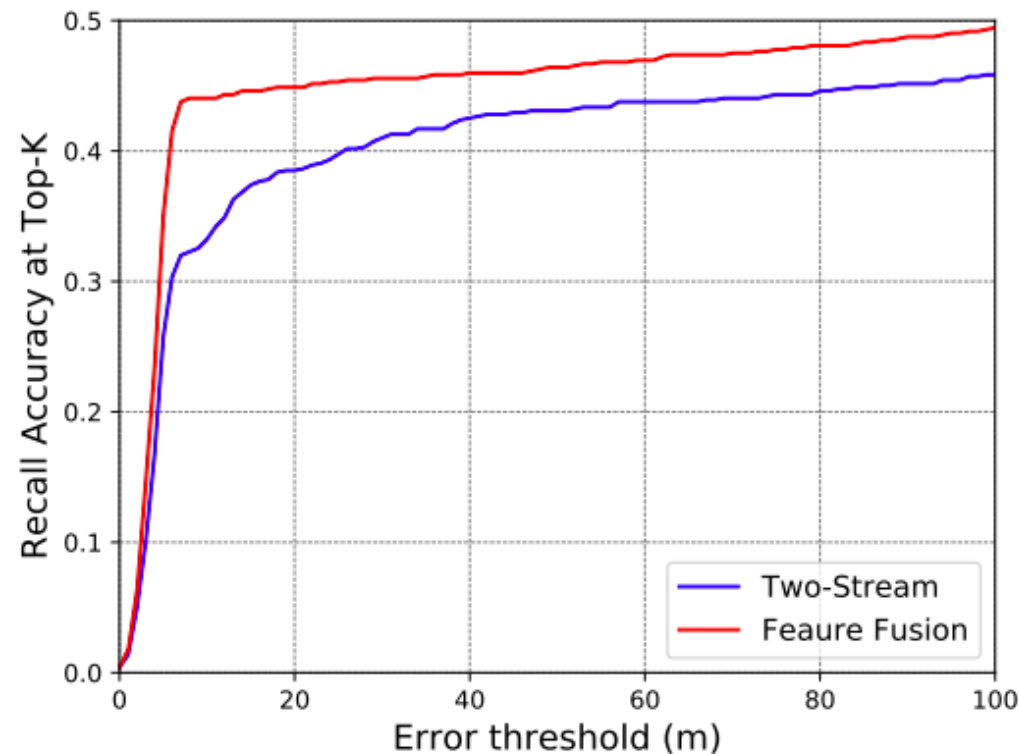


Image Retrieval (OP Dataset)

Ground Query

Top matches (top 1 – top 5 from left to right)



0.0 m



11.08 m



44.33 m



254.90 m



69.32 m



2.25 m



161.17 m



12.35 m



21.19 m



521.67 m



4.44 m



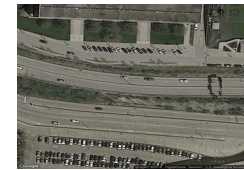
111.39 m



9.58 m



1246.39 m



424.53 m



4.75 m



191.66 m



154.27 m



256.27 m



81.91 m

Aerial-to-Ground (A2G) Image Matching

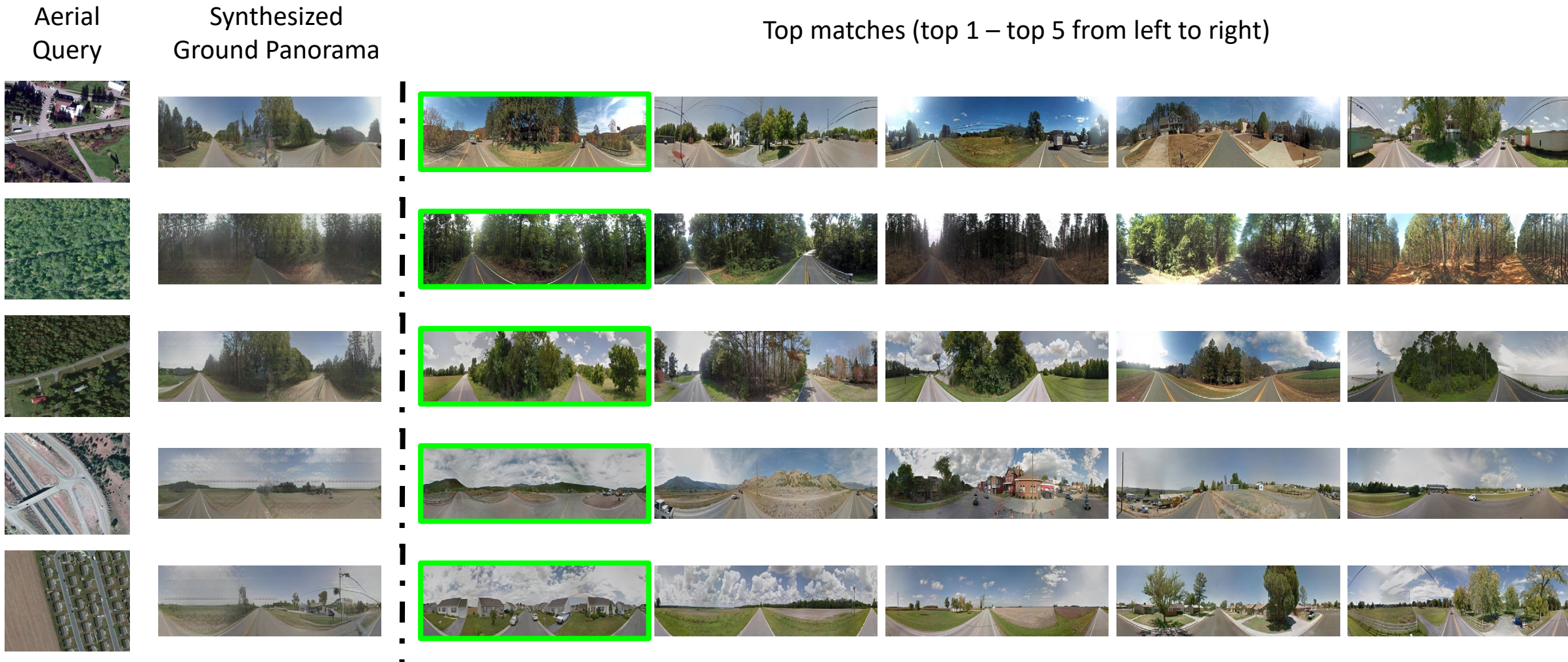
Reverse problem

Synthesize the ground panorama from aerial image

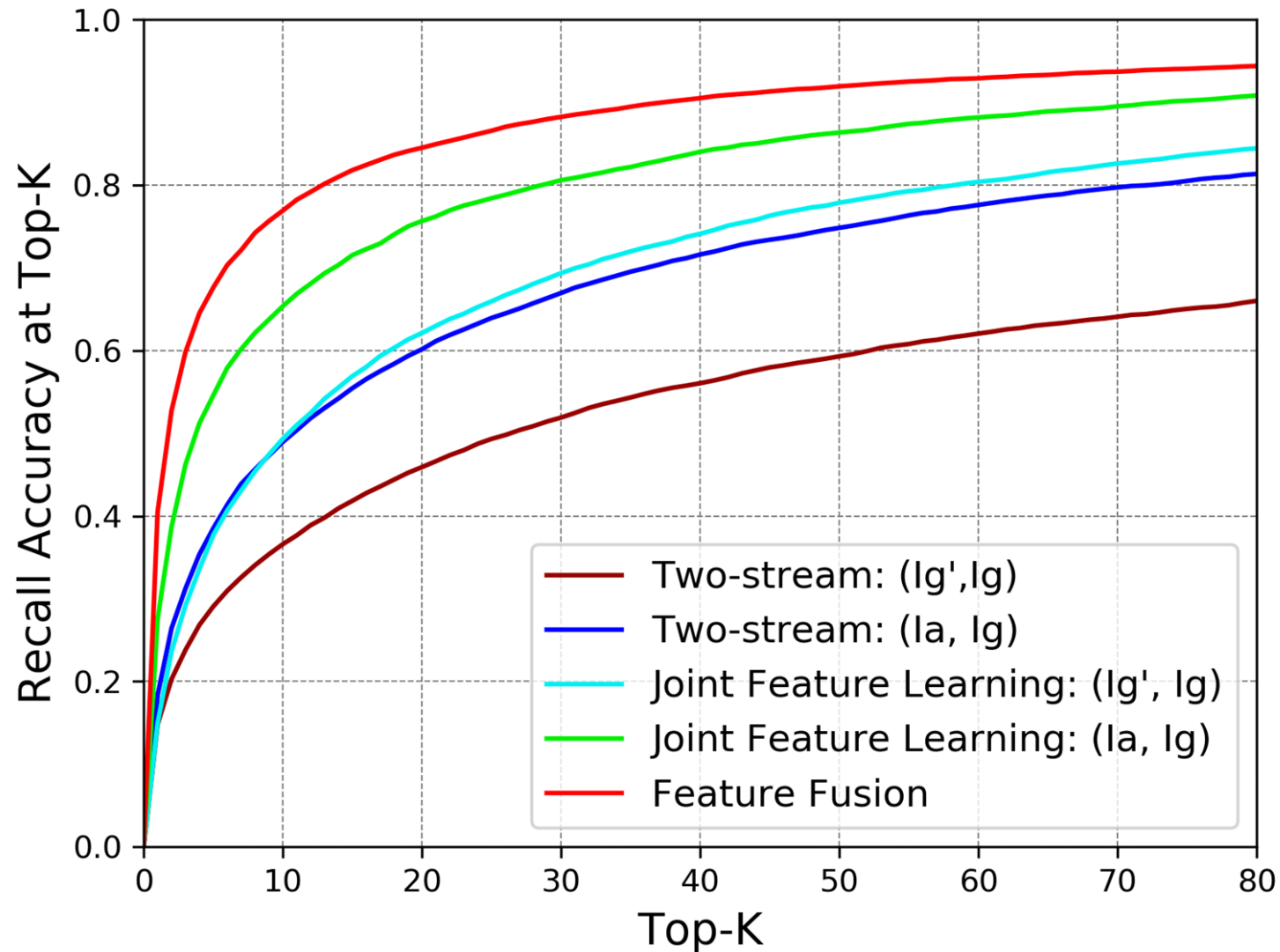
Quantitative Results - A2G (CVUSA Dataset)

Method	Top-1	Top-10	Top-1%
Two-stream baseline ($I_{g'}$, I_g)	15.04%	37.31%	67.99%
Two-stream baseline (I_a , I_g)	16.99%	47.06%	82.11%
Joint Feat. Learning ($I_{g'}$, I_g)	16.46%	50.26%	86.26%
Joint Feat. Learning (I_a , I_g)	27.39%	65.29%	91.46%
Feature Fusion	44.99%	79.37%	95.66%

Image Retrieval - A2G (CVUSA Dataset)



Recall Accuracy - A2G (CVUSA Dataset)



Summary

- Novel and practical approach to cross-view image matching
- Domain gap is bridged by synthesized images
- Significant improvement on Top-1 and Top-10 accuracies over SOTA on CVUSA.
- This approach can be used for other view transformation tasks where the transformations can be in horizontal or vertical directions.



CENTER FOR RESEARCH
IN COMPUTER VISION



ICCV 2019
Seoul, Korea

Bridging the Domain Gap for Ground-to-Aerial Image Matching



Krishna Regmi

Krishna Regmi & Mubarak Shah
University of Central Florida

Summary

- Pixel-Wise Geo-localization
 - Geodetic Alignment of Aerial Video Frames
- Image-Based Geo-Localization
 - Same View (Street-View to Street-View)
 - Generalized Maximum Clique (PAMI, 2014)
 - Constraint Dominant Sets (PAMI, 2017)
 - Cross-View Geo-Localization
 - Bird's Eye-View to Street View (CVPR, 2017)
 - Aerial to Ground View (ICCV, 2019)

Thank You