

# Vehicle localization without SLAM: Learning to find your camera's pose in an aerial image

Julian Kooij

Intelligent Vehicles group, TU Delft

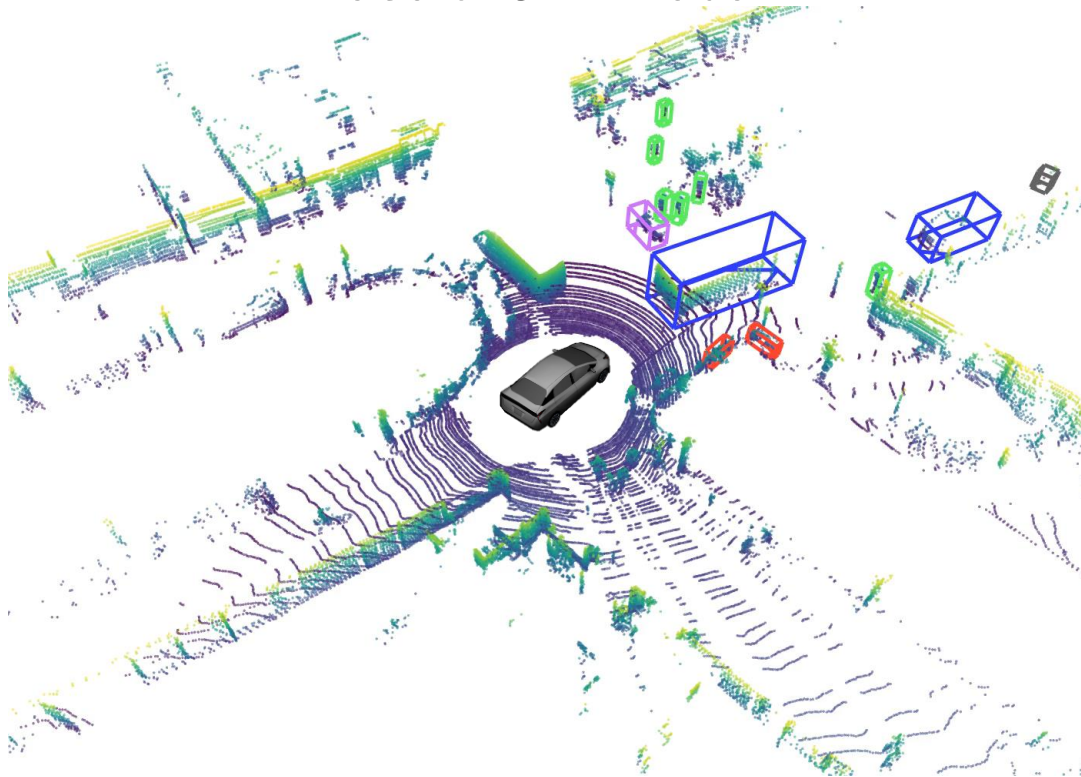
# Julian Kooij

- **Associate Professor**  
Intelligent Vehicles group  
Cognitive Robotics department  
Mechanical Engineering Faculty, TU Delft
- **Director of 3DUU Delft AI lab**  
TUD AI Initiative Labs & Talent program  
Collaboration between ME and ABE faculties
- **PhD**  
*2010-2014*  
University of Amsterdam  
Daimler R&D (Mercedes), Ulm Germany



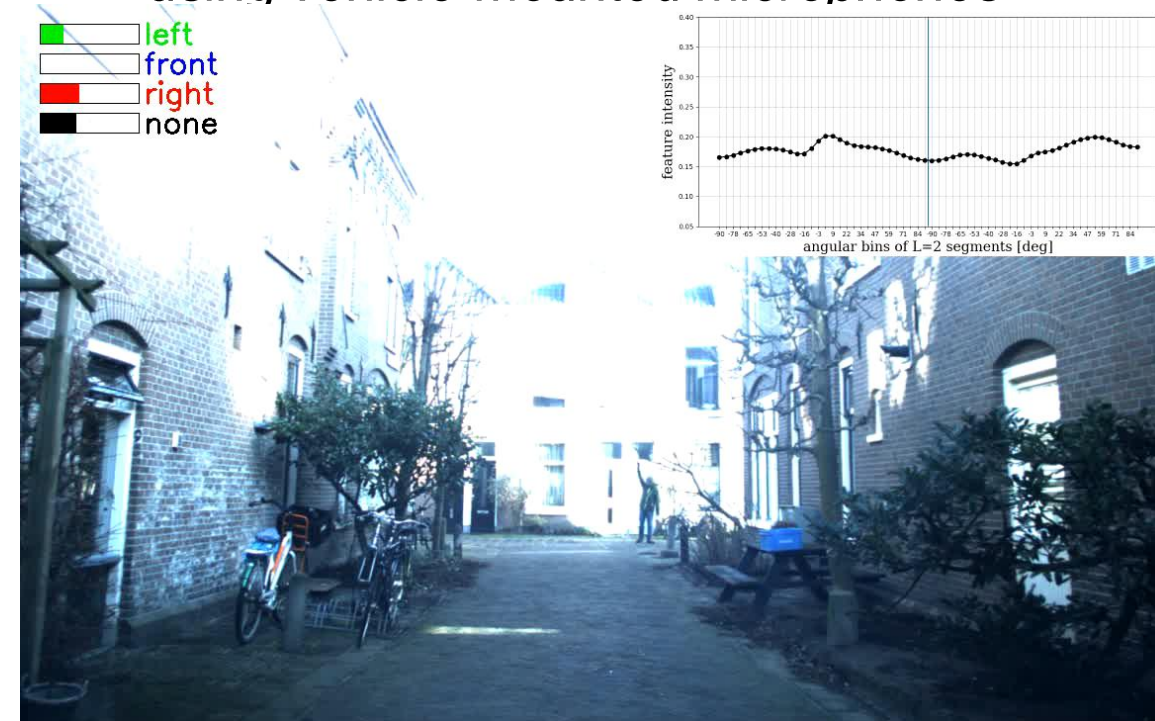
# Research at IV group

*Driving around Vulnerable Road Users,  
use of 3+1D radar*



*Multi-class Road User Detection with 3+1D Radar in the View-of-Delft Dataset, A. Palffy, E.A.I. Pool, S. Baratam, J.F.P. Kooij, D.M. Gavrila, IEEE Robotics and Automation Letters, 2022, vol. 7(2), 4961-4968*

*Predict approaching traffic  
using vehicle-mounted microphones*

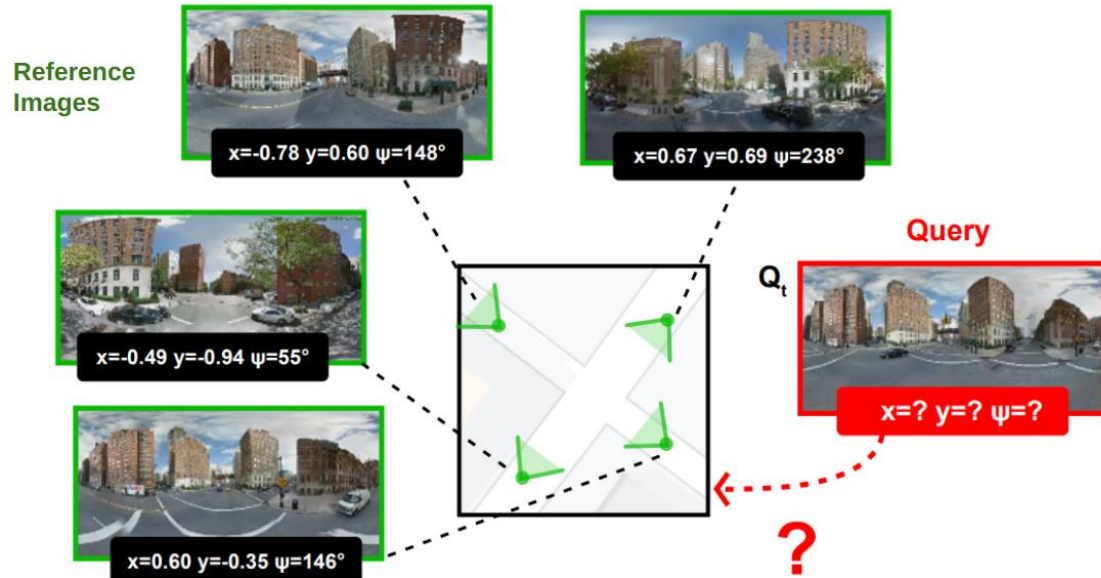


*Hearing What You Cannot See: Acoustic Vehicle Detection Around Corners, Y. Schulz, A.K. Mattar, T.M. Hehn, J.F.P. Kooij, IEEE Robotics and Automation Letters (RA-L), 2021, vol. 6(2), 2587-2594*



# 3D Urban Understanding (3DUU) AI lab

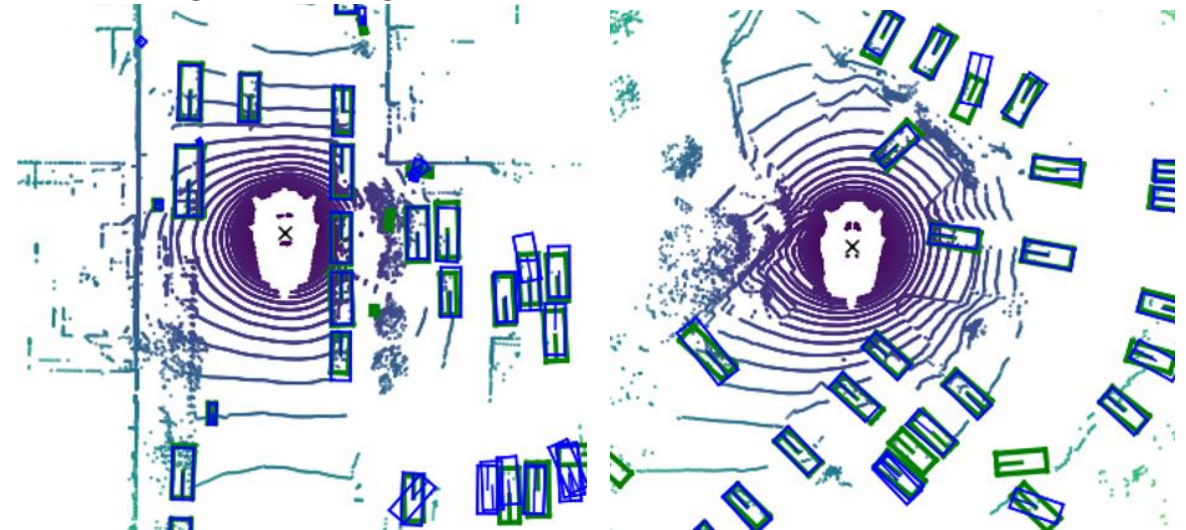
*Localization: Where was this 2D image taken in 3D world?*



Reconstructing & imagining



*Fusing 2D images and 3D point clouds for detection*



*Interpreting and segmenting pointclouds*





# TU Delft Intelligent Vehicles group

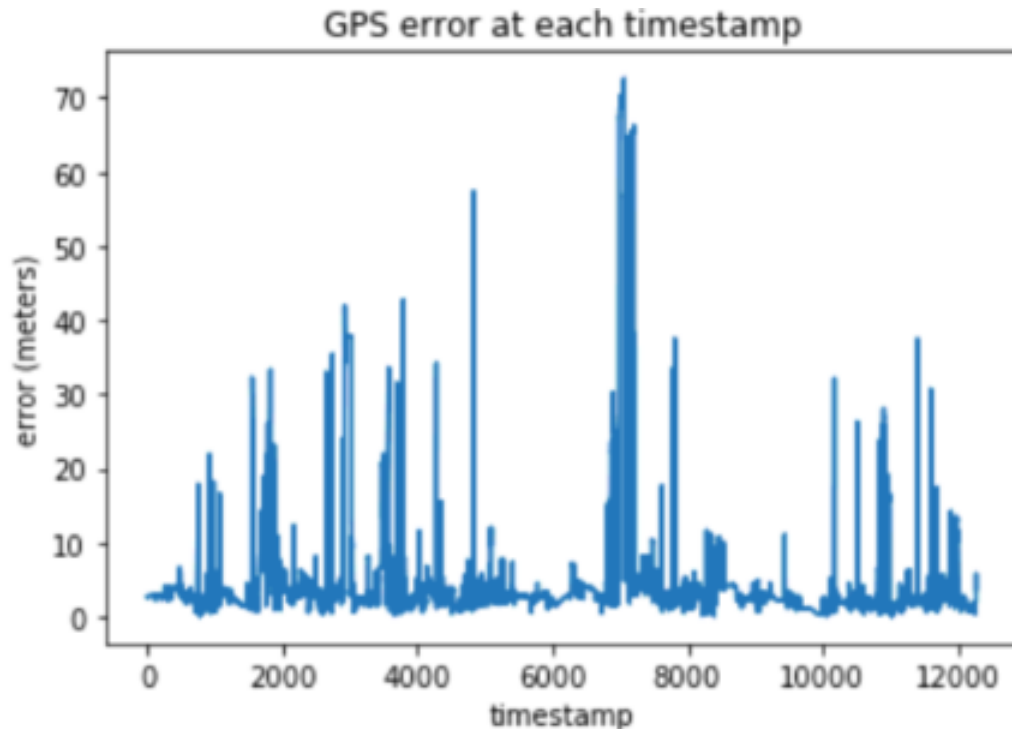
<https://intelligent-vehicles.org/>



# The problem of GPS

- GNSS/GPS → can reach tens of meters<sup>1</sup> in urban areas

1. Ben-Moshe, Boaz, et al. "Improving Accuracy of GNSS Devices in Urban Canyons." CCCG. 2011



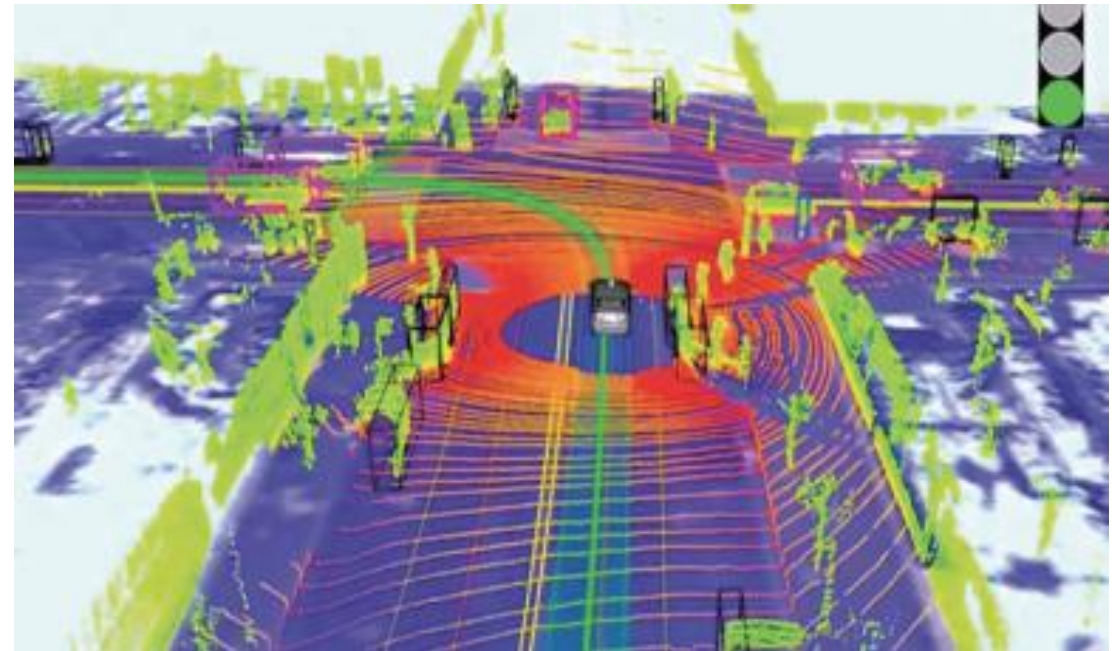
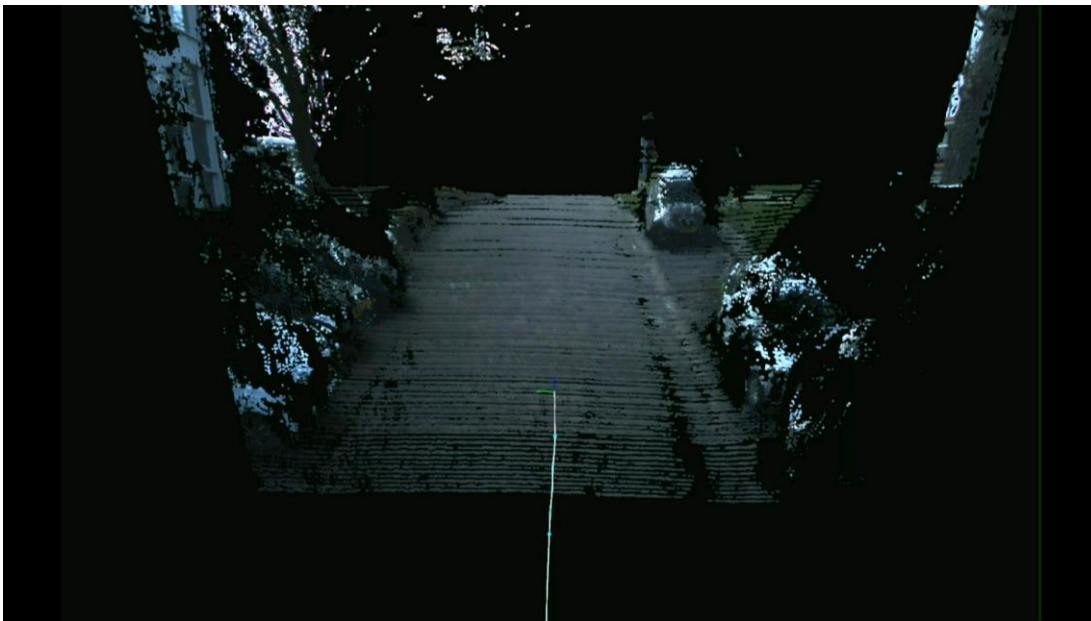
One traversal from Oxford RobotCar

- Red: Raw GPS
- Green: Ground truth RTK  
(high precision GPS, post-processing)
- Here, worst-case ~70 meters



## Model-driven visual localization

- Alternative: match sensor data to detailed 3D models,  
but implies huge data collection, processing, updating effort ... hard to scale



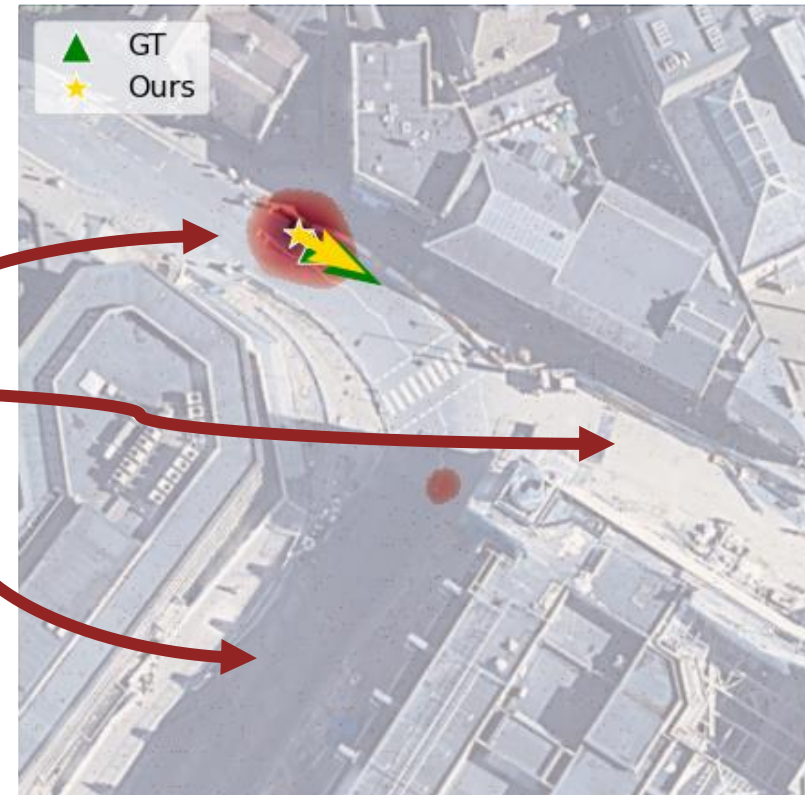
Credit right image: "Autonomous driving in the iCity—HD maps as a key challenge of the automotive industry.", Seif, Heiko G., and Xiaolong Hu, *Engineering* 2.2 (2016): 159-162.

## Data-driven visual localization

- Can we not use aerial images (“Google Earth”) as easily obtained map?
- Match learned representations of vehicle’s camera image to map’s aerial image



Does it  
match?



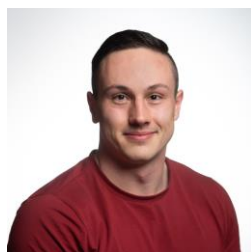


JUNE 18-22, 2023

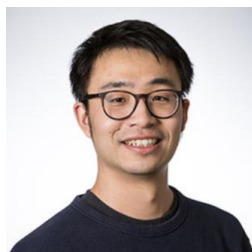
CVPR



# SliceMatch: Geometry-guided Aggregation for Cross-View Pose Estimation CVPR 2023



Ted Lentsch\*



Zimin Xia\*



Holger Caesar

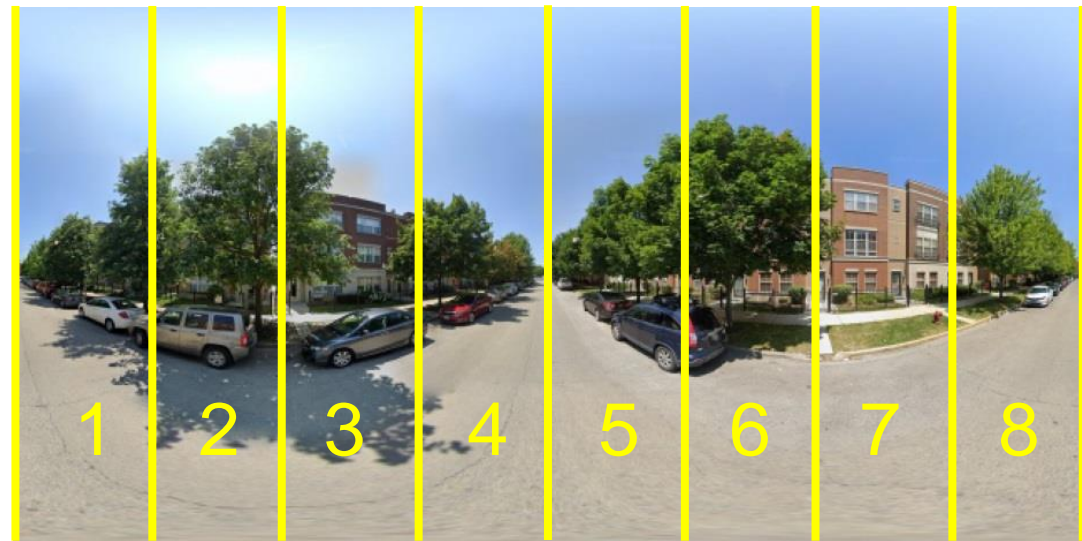


Julian F. P. Kooij

Intelligent Vehicles Group, Delft University of Technology, The Netherlands

\*Equal contribution

## What is a slice?

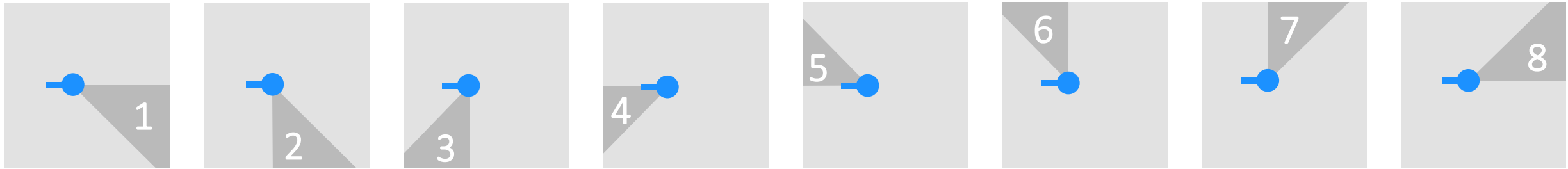


Cosine Similarity



# Geometry-guided feature aggregation

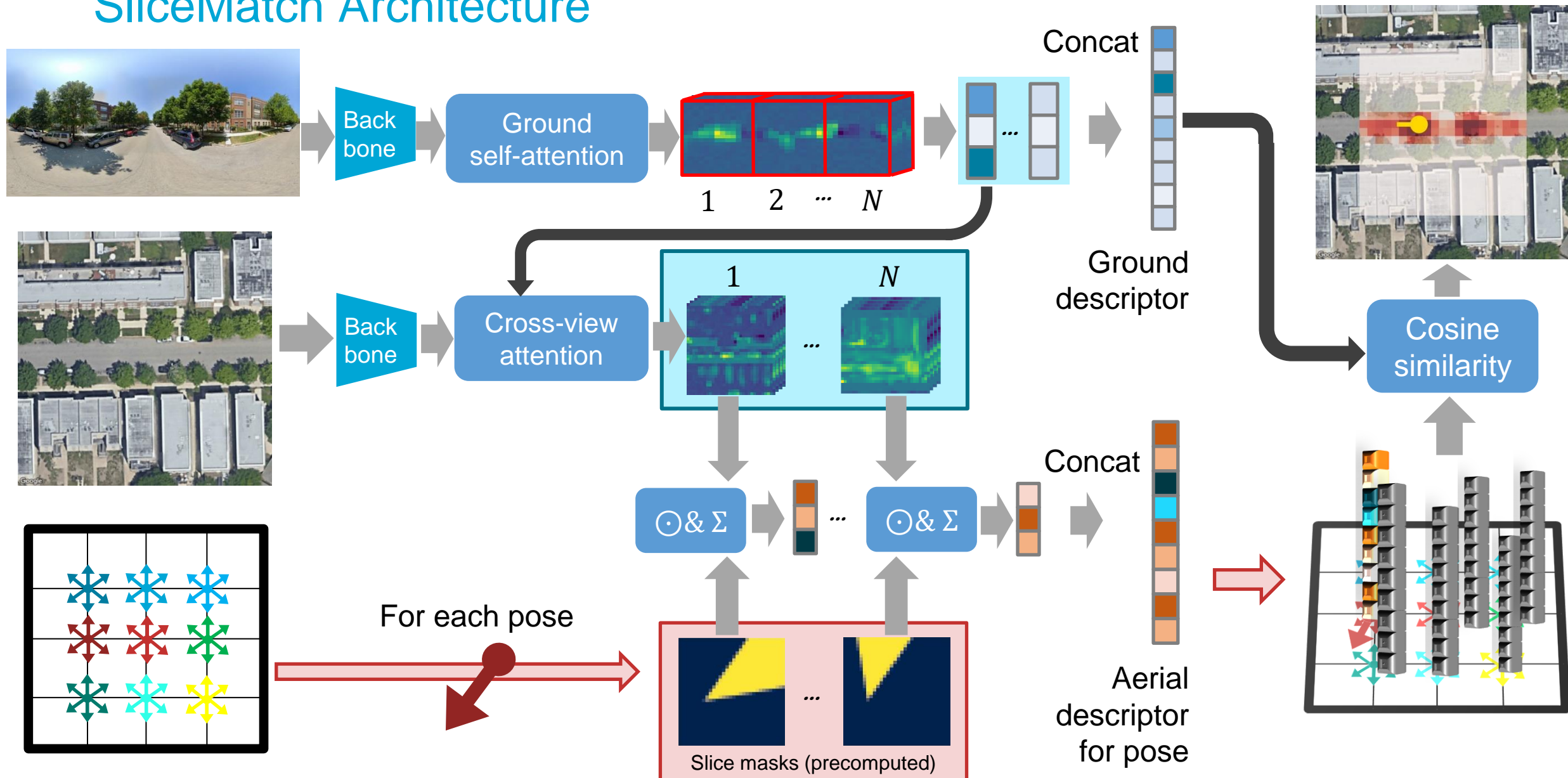
- **Goal:** Achieve accuracy by testing many pose hypotheses
- **Idea:** Minimize computations per pose hypothesis → only do feature aggregation per pose
- **“Slice masks”:** Per pose, pre-compute the aerial region for each slice



- Use contrastive learning to obtain pos + orientation discriminative features



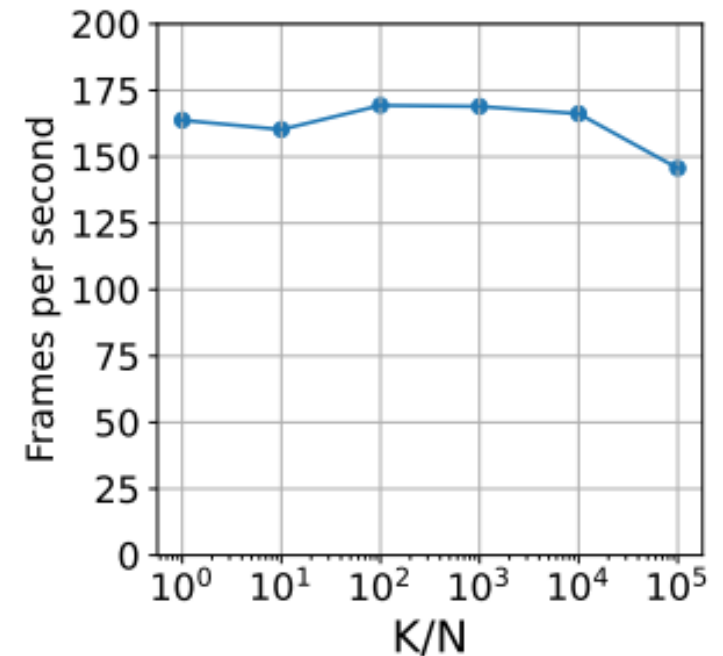
# SliceMatch Architecture





## Efficient computation

- Feature extraction is done only once
- Per pose hypotheses operations are minimized:
  - Feature aggregation module uses pre-computed slice masks
  - Apply slice masks to feature maps: just multiply and sum
  - Highly parallelizable tensor operations
- Inference speed at over 150 FPS (on V100 GPU)
- Significantly faster than previous SOTA method
- Did not observe run time increase when testing more poses



$K$  = number of poses  
 $N$  = 16 slices

## Quantitative evaluation

Model	Backbone	Aligned Images	Same-Area				Cross-Area			
			↓ Location (m)		↓ Orientation (°)		↓ Location (m)		↓ Orientation (°)	
			Mean	Median	Mean	Median	Mean	Median	Mean	Median
CVR [55]	VGG16	✓	8.99	7.81	-	-	8.89	7.73	-	-
MCC [50]	VGG16	✓	6.94	3.64	-	-	9.05	5.14	-	-
SliceMatch (ours)	VGG16	✓	<b>5.18</b>	<b>2.58</b>	-	-	<b>5.53</b>	<b>2.55</b>	-	-
MCC [50]	VGG16	X	9.87	6.25	56.86	16.02	12.66	9.55	72.13	29.97
SliceMatch (ours)	VGG16	X	<b>8.41</b>	<b>5.07</b>	<b>28.43</b>	<b>5.15</b>	<b>8.48</b>	<b>5.64</b>	<b>26.20</b>	<b>5.18</b>
SliceMatch (ours)	ResNet50	X	<b>6.49</b>	<b>3.13</b>	<b>25.46</b>	<b>4.71</b>	<b>7.22</b>	<b>3.31</b>	<b>25.97</b>	<b>4.51</b>

CVR: Zhu, et al, CVPR 2021

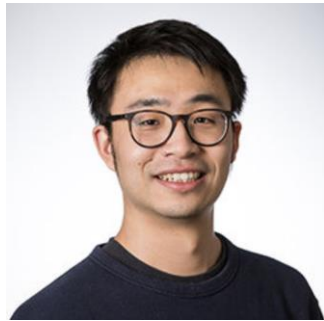
MCC: Xia, et al, ECCV 2022

- SliceMatch outperforms SOTA baselines (at the time) in all use cases
- With a more advanced backbone (ResNet50 vs VGG16), accuracy further improves
- *(N.B.: we improved the GT localization annotations for the VIGOR datasets)*



# Convolutional Cross-View Pose Estimation

IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)



Zimin Xia



Olaf Booij



Julian F. P. Kooij



Intelligent Vehicles Group, TU Delft, The Netherlands

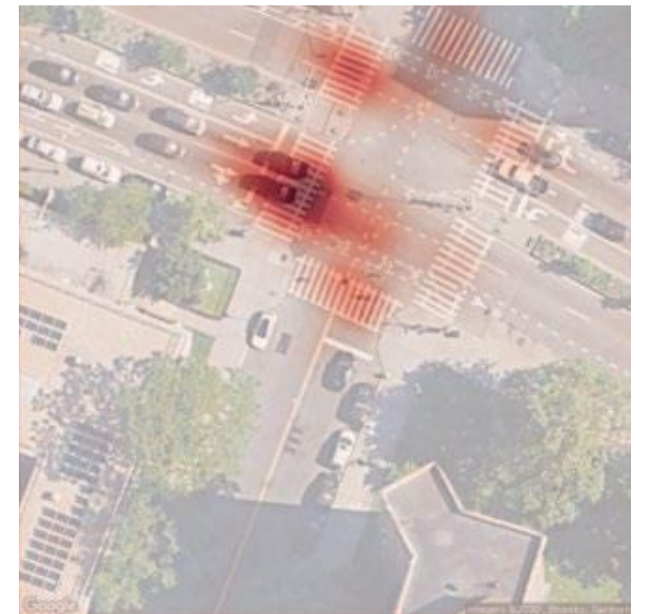
TomTom, Amsterdam, The Netherlands

# Cross-view localization

- Pinpoint the location of the ground camera on a local aerial image
- We formulate localization as dense multi-class classification

 $G$  $A$ 

$$p(D|G, A)$$

 $D$

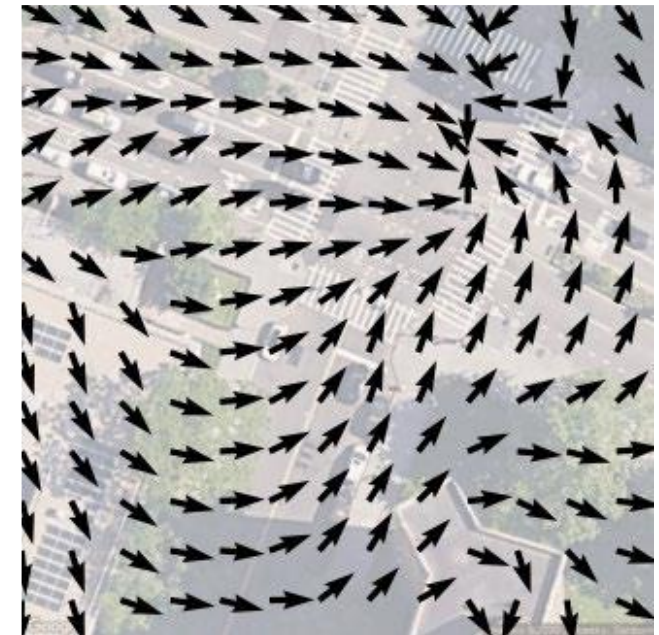


# Cross-view pose estimation

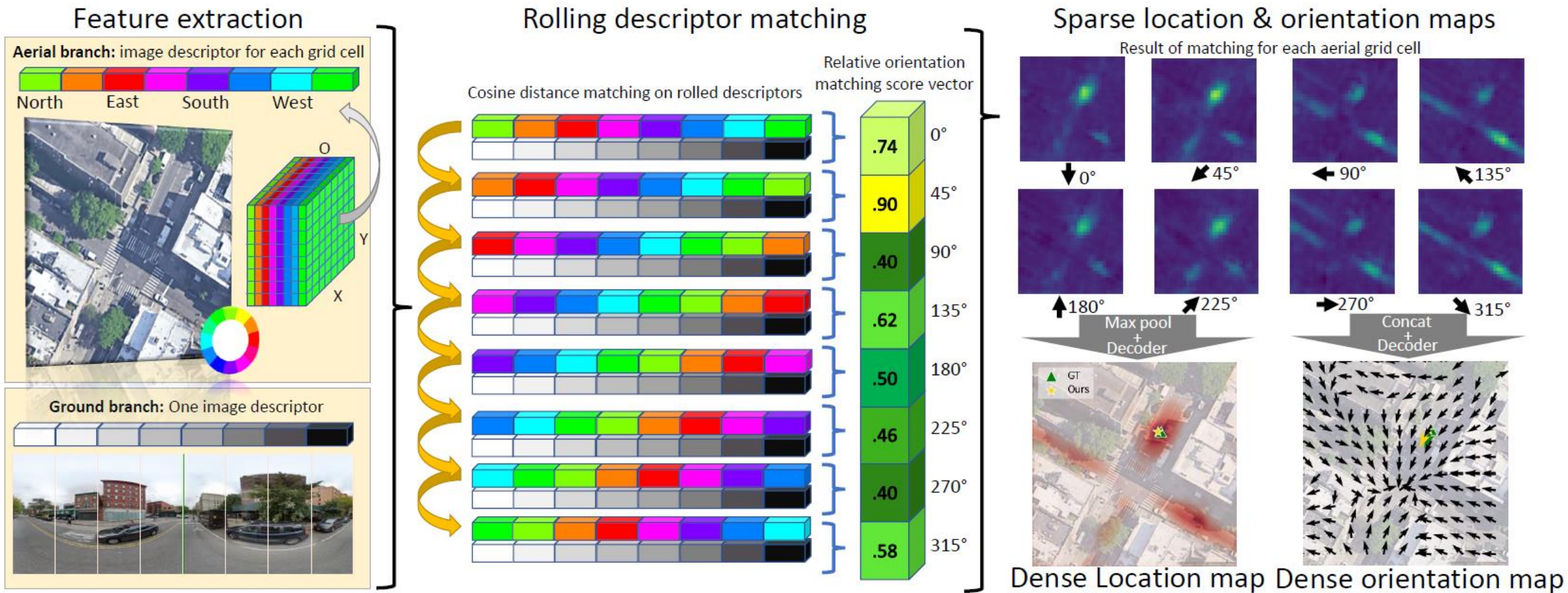
- Estimate location and *orientation*
- Location and orientation are related
- Solution: also predict a vector field

 $G$  $A$ 

$$f(G, A) \rightarrow O$$

 $O$

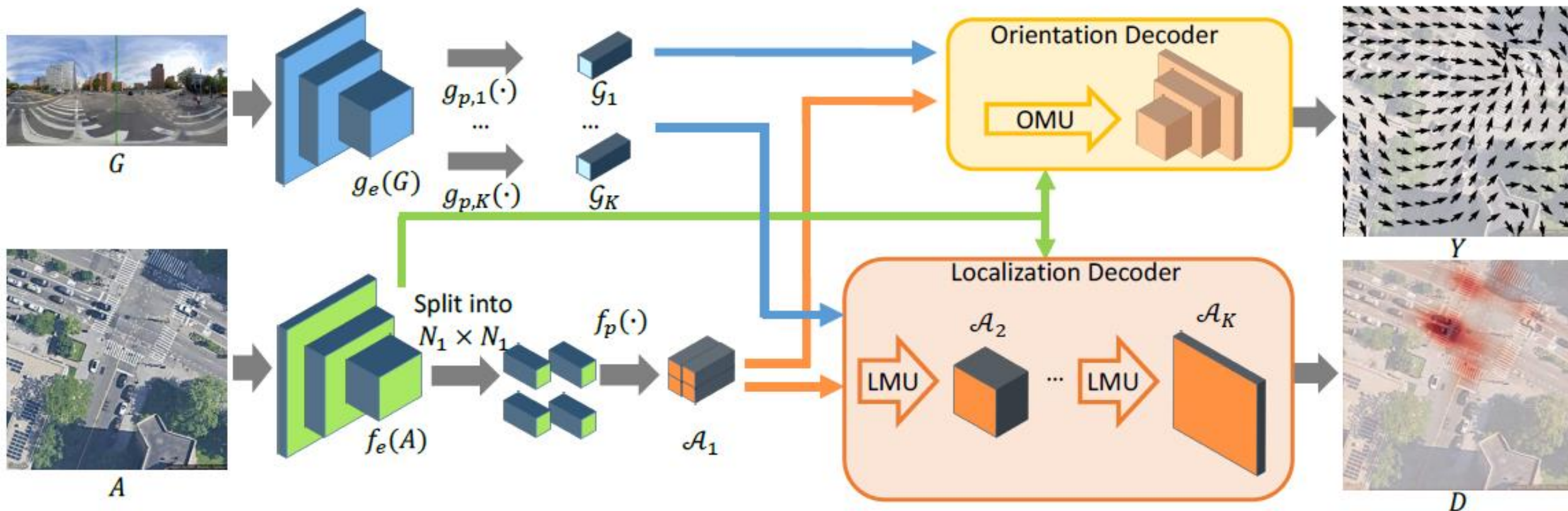
# Architecture





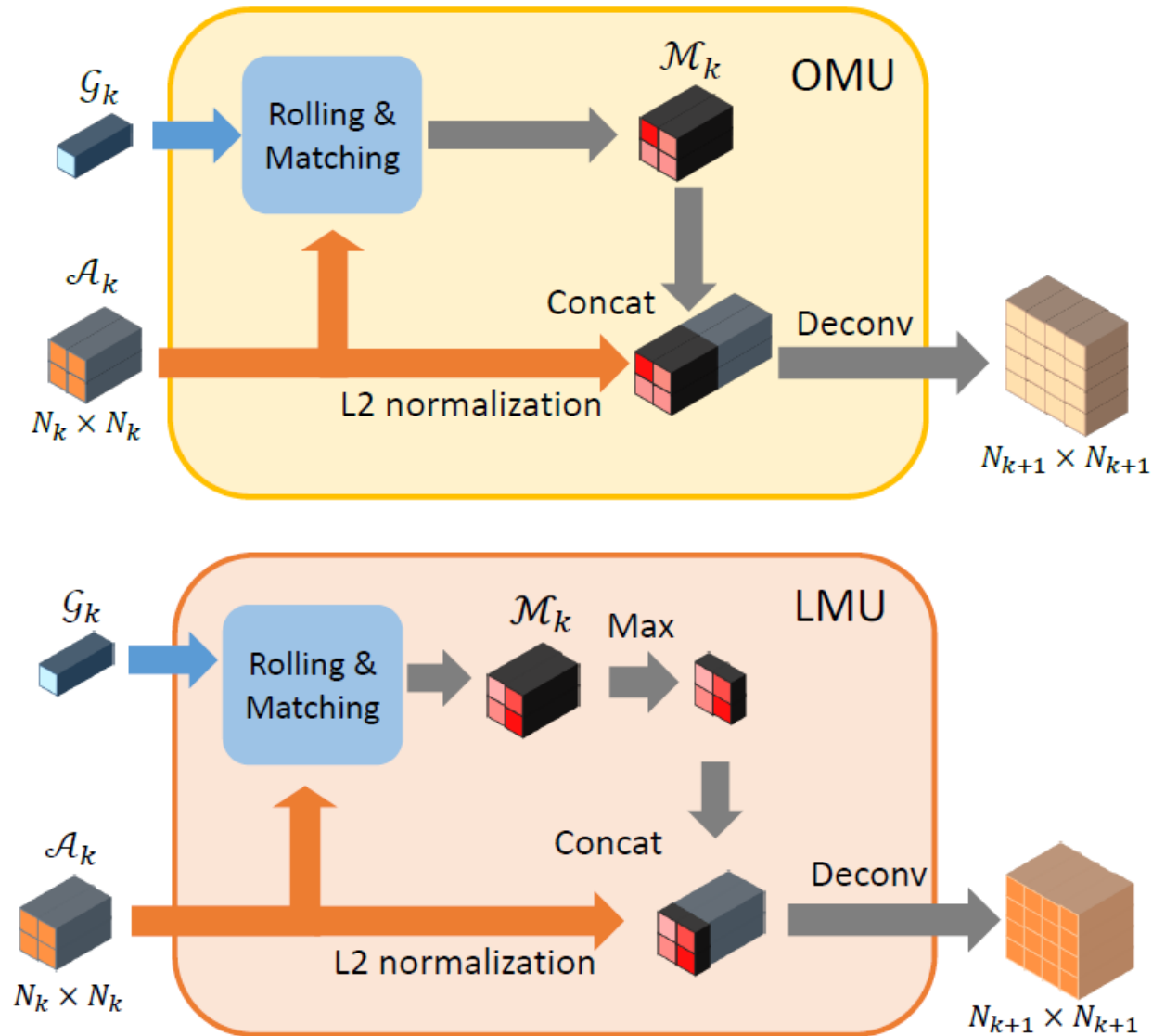
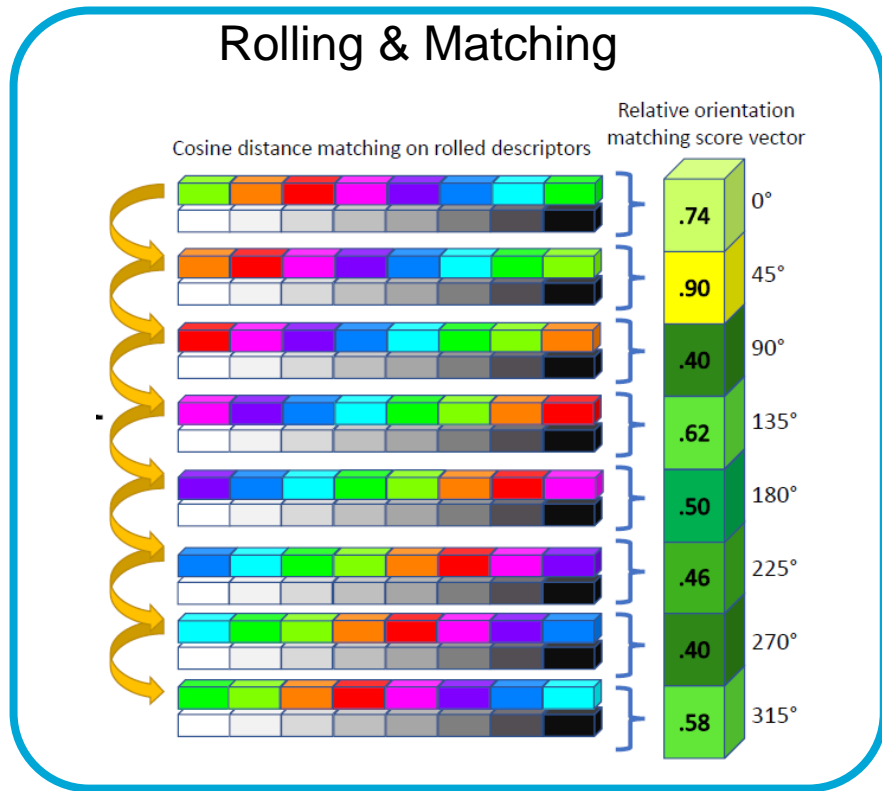
# Architecture

- Siamese-like network, with two predictor heads
- Rolling & Matching happens in the Decoders



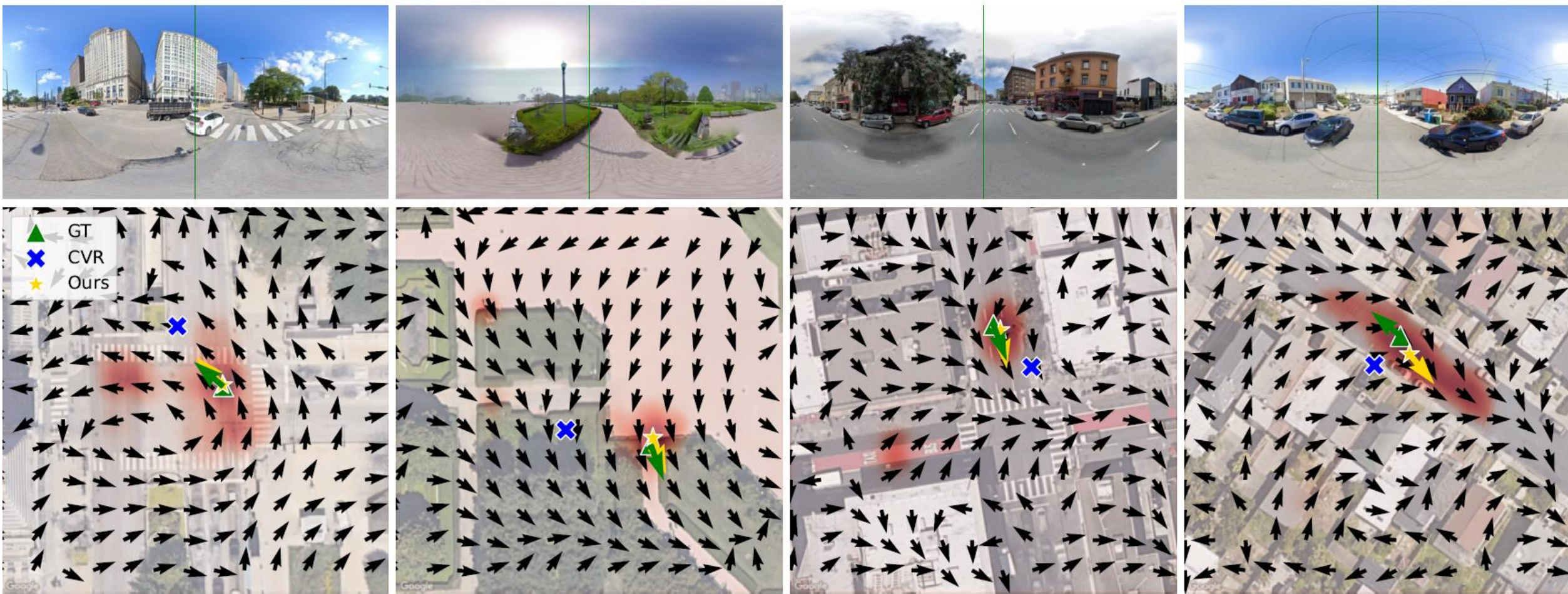


# Architecture



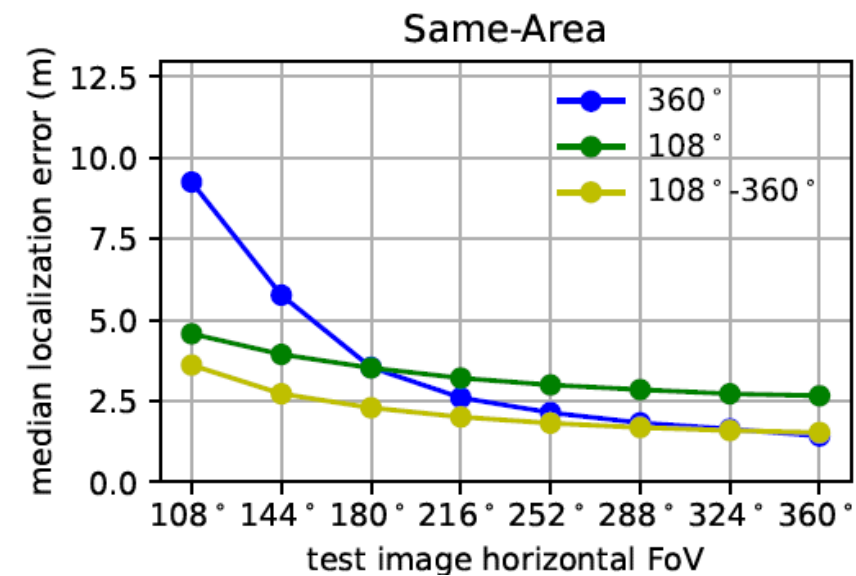
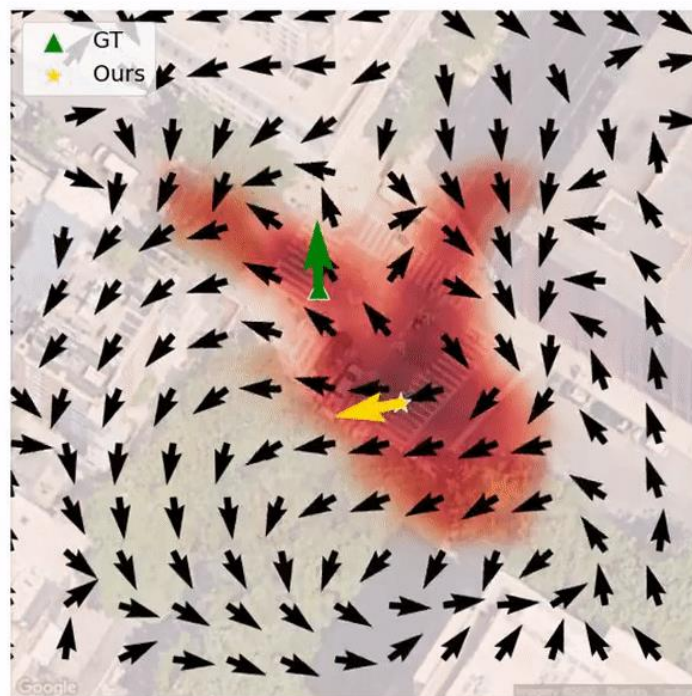
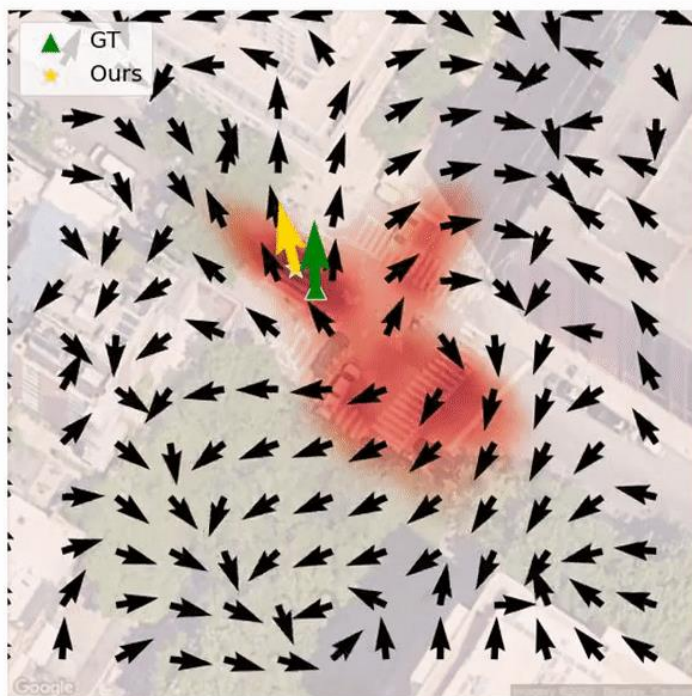
- OMU = Orientation Matching Upsampling
- LMU = Localization Matching Upsampling

# Qualitative results



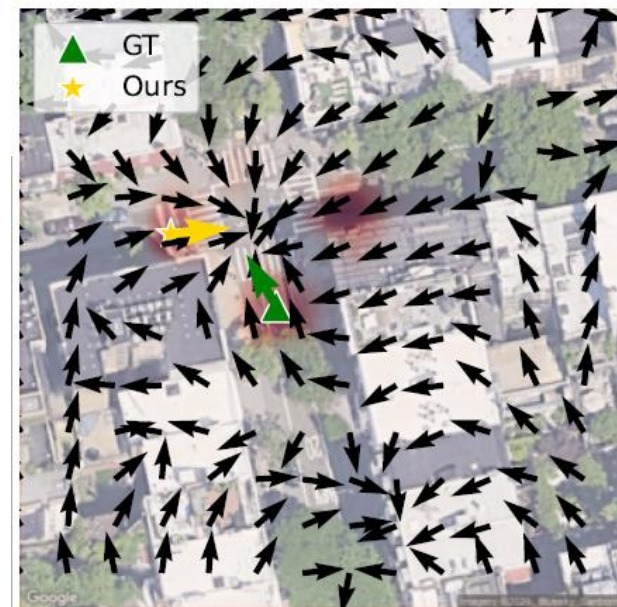


## Qualitative results, varying FoV

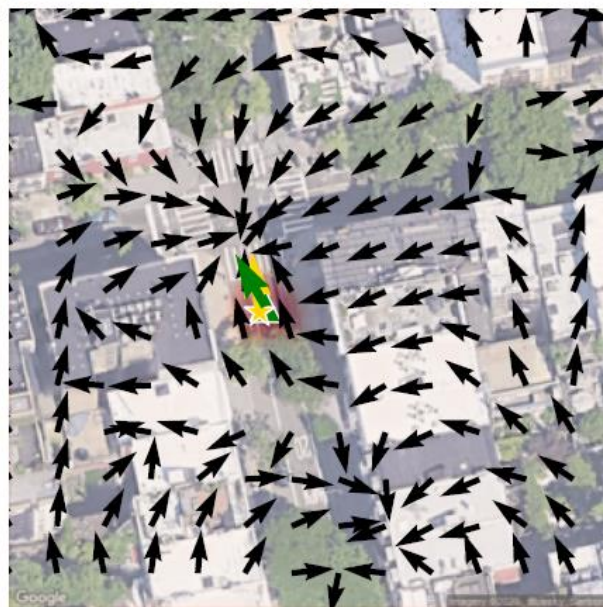




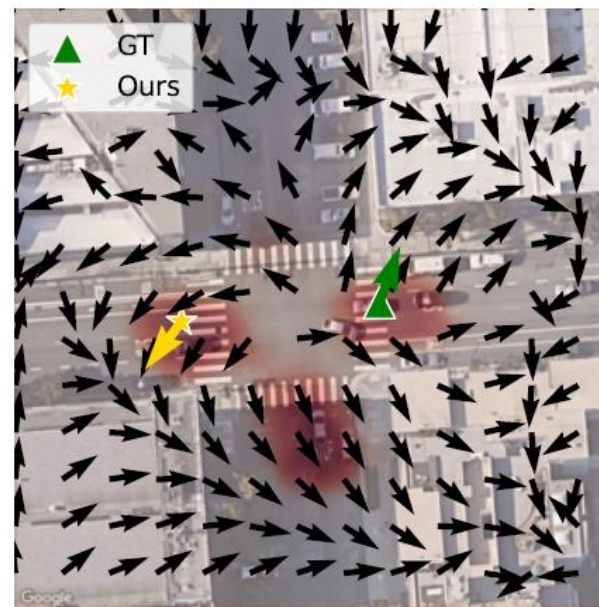
# Adding orientation prior



no orientation prior



known orientation



no orientation prior



known orientation



# Quantitative results

- VIGOR**

Large scale dataset, 4 cities in USA

Same-area: train 2 cities, test same cities

Cross-area: train 2 cities, test other cities

VIGOR test		Same-Area						Cross-Area					
		↓ Localization (m)		↓ Orientation (°)		↑ P@GT ( $10^{-3}$ )		↓ Localization (m)		↓ Orientation (°)		↑ P@GT ( $10^{-3}$ )	
		mean	median	mean	median	mean	median	mean	median	mean	median	mean	median
0°	CVR [22]	8.82	7.68	-	-	0.02	0.02	9.45	8.33	-	-	0.02	0.02
	Eff-CVR	7.89	6.25	-	-	0.02	0.03	8.27	6.60	-	-	0.02	0.03
	SliceMatch [32]	5.18	2.58	-	-	0.06	0.05	5.53	2.55	-	-	0.06	0.06
	CCVPE (ours)	<b>3.60</b>	<b>1.36</b>	-	-	<b>1.60</b>	<b>1.12</b>	<b>4.97</b>	<b>1.68</b>	-	-	<b>1.08</b>	<b>0.71</b>
360°	SliceMatch [32]	8.41	5.07	28.43	<b>5.15</b>	0.02	0.02	8.48	5.64	<b>26.20</b>	<b>5.18</b>	0.02	0.02
	CCVPE (ours)	<b>3.74</b>	<b>1.42</b>	<b>12.83</b>	6.62	<b>1.47</b>	<b>1.00</b>	<b>5.41</b>	<b>1.89</b>	<b>27.78</b>	13.58	<b>0.93</b>	<b>0.58</b>

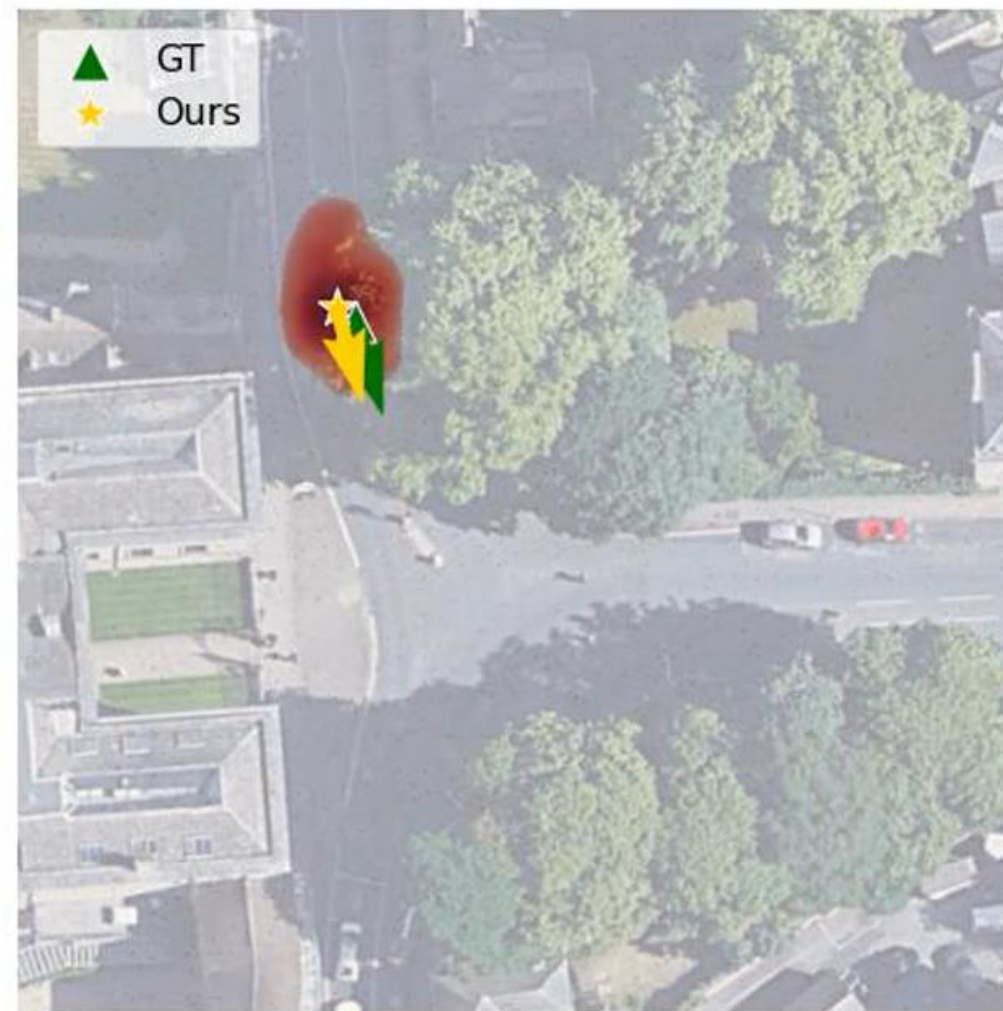
- KITTI**

Same-Area		↓ Localization (m)		↑ Lateral (%)			↑ Longitudinal (%)			↓ Orientation (°)		↑ Orientation (%)		
		mean	median	1m	3m	5m	1m	3m	5m	mean	median	1°	3°	5°
retrieval	CVM-Net [16]	-	-	5.83	17.41	28.78	3.47	11.18	18.42	-	-	-	-	-
	CVFT [19]	-	-	7.71	22.37	36.28	3.82	11.48	18.63	-	-	-	-	-
	SAFA [15]	-	-	9.49	29.31	46.44	4.35	12.46	21.10	-	-	-	-	-
	Polar-SAFA [15]	-	-	9.57	30.08	45.83	4.56	13.01	21.12	-	-	-	-	-
	DSM [47]	-	-	10.12	30.67	48.24	4.08	12.01	20.14	-	-	3.58	13.81	24.44
±10°	LM [39]	12.08	11.42	35.54	70.77	80.36	5.22	15.88	26.13	3.72	2.83	19.64	51.76	71.72
	SliceMatch [32]	7.96	4.39	49.09	91.76	98.52	15.19	49.99	57.35	4.12	3.65	13.41	42.62	64.17
	CCVPE (ours)	<b>1.22</b>	<b>0.62</b>	<b>97.35</b>	<b>98.65</b>	<b>99.71</b>	<b>77.13</b>	<b>96.08</b>	<b>97.16</b>	<b>0.67</b>	<b>0.54</b>	<b>77.39</b>	<b>99.47</b>	<b>99.95</b>
360°	LM [39]	15.51	15.97	5.17	15.13	25.44	4.66	15.00	25.39	89.91	90.75	0.61	1.88	2.89
	SliceMatch [32]	9.39	5.41	39.73	<b>80.56</b>	<b>87.92</b>	13.63	40.75	49.22	<b>8.71</b>	<b>4.42</b>	<b>11.35</b>	<b>36.23</b>	<b>55.82</b>
	CCVPE (ours)	<b>6.88</b>	<b>3.47</b>	<b>53.30</b>	77.63	85.13	<b>25.84</b>	<b>55.05</b>	<b>68.49</b>	15.01	6.12	8.96	26.48	42.75

# Cross-view pose estimation on Oxford RobotCar



*N.B.: Single frame results,  
no temporal filtering, no sensor fusion!*





# Conclusions

Cross-view pose estimation is a rapidly developing field

- Current focus is on fine-grained localization, relying on rough location prior
- Originally only localization, now also orientation estimation

Use Multi-modal distributions

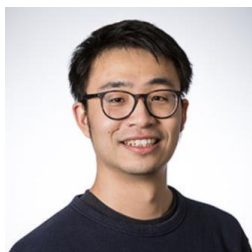
- Represent localization ambiguity, supports confidence estimation
  - Can condition orientation on localization estimate
- 
- On Oxford RobotCar, we start to reach  $< 1\text{m}$  performance. VIGOR dataset is more difficult.
  - Evaluation is single frame only, temporal and sensor fusion should highly improve accuracy



EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO  
2024

# Adapting Fine-Grained Cross-View Localization to Areas without Fine Ground Truth



Zimin Xia<sup>1</sup>



Yujiao Shi<sup>2</sup>



Hongdong Li<sup>3</sup>



Julian F. P. Kooij<sup>4</sup>

École Polytechnique Fédérale de Lausanne (EPFL)

ShanghaiTech University

Australian National University

Delft University of Technology

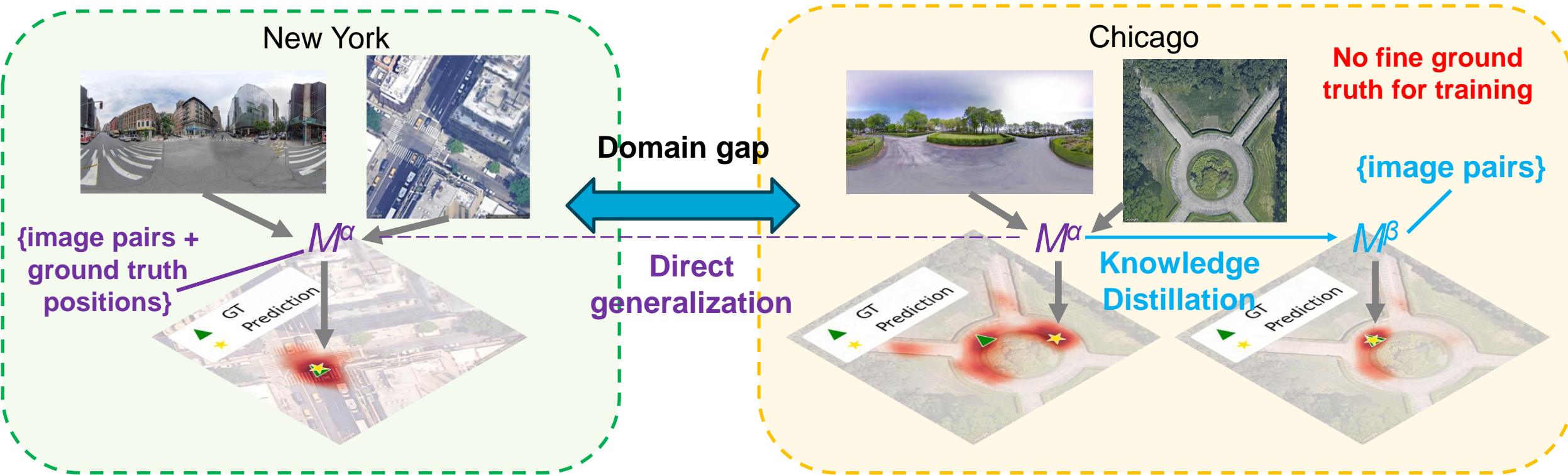


上海科技大学  
ShanghaiTech University



Australian  
National  
University

# Does CV pose estimation generalize to new unseen areas?



- Collecting fine ground truth in a new target city can be expensive or infeasible
- Collecting images with noisy location data is easy, e.g. by phone-grade GPS

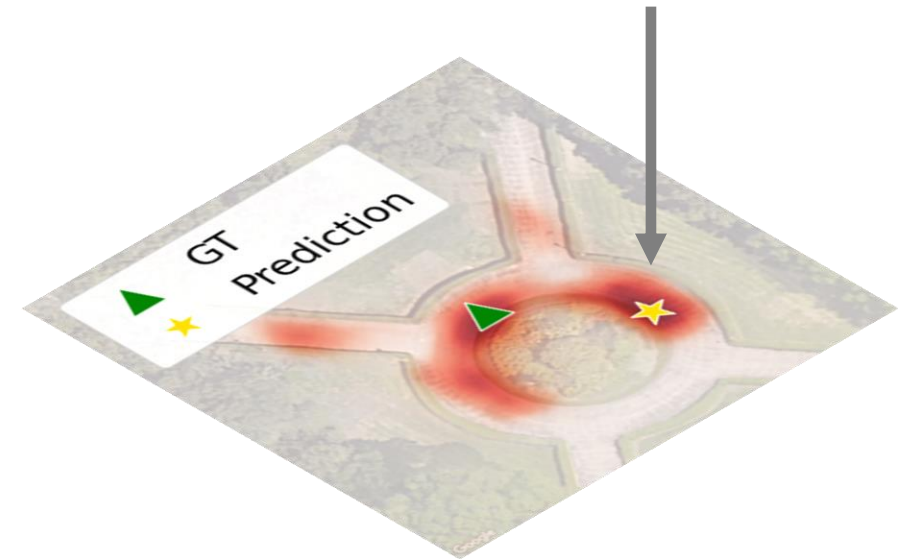


# Leveraging noisy location data to adapt the models to new areas

## Observations:

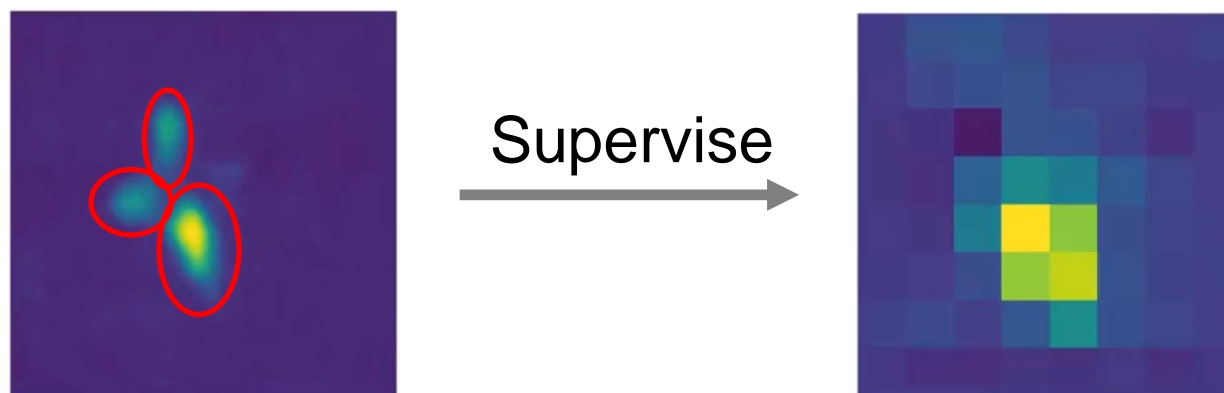
- Direct generalization leads to uncertain localization:
  - high uncertainty
  - small positional noise
  - large outliers
- SOTA<sup>[1,2]</sup> predicts heat maps for localization
- SOTA<sup>[1,2]</sup> has coarse-to-fine outputs

Use these insights to design a suitable  
knowledge distillation strategy



# Leveraging noisy location data to adapt the models to new areas

- Use the pretrained model as a teacher to supervise a copy (student) of itself
- But how can this really improve the model?
  1. Use teacher's high-res output to supervise student's lower-resolution heat maps (remember: SotA models use coarse-to-fine strategy)
  2. Take teacher's multi-modal output and select single mode as "pseudo GT"  
*"pretend the teacher's best guess is correct"*

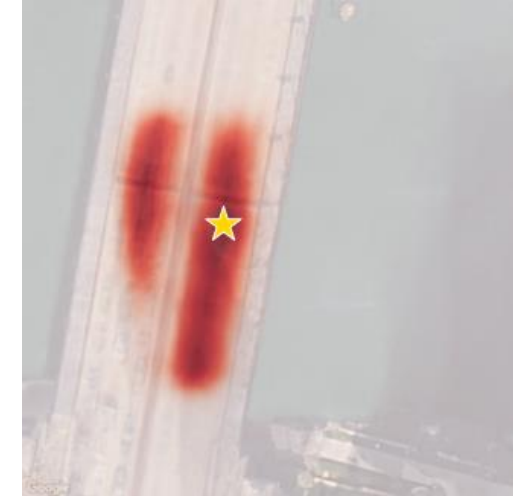


# Leveraging noisy location data to adapt the models to new areas

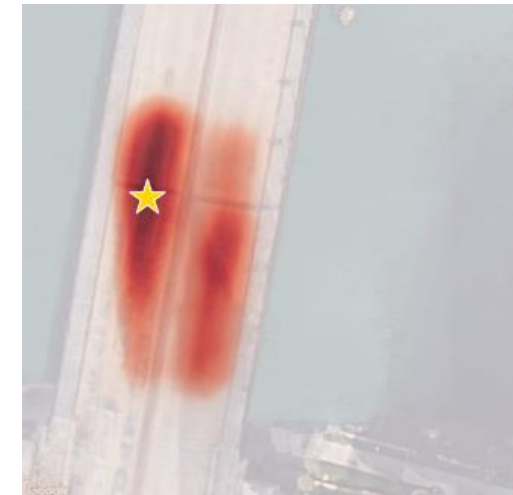
What if there are multiple equally likely modes?

Train an auxiliary student model for outlier detection:

- Measure difference between the teacher's and the student's predictions
- If the difference is large, exclude this sample in final student model training
- Keep top X% more consistent samples



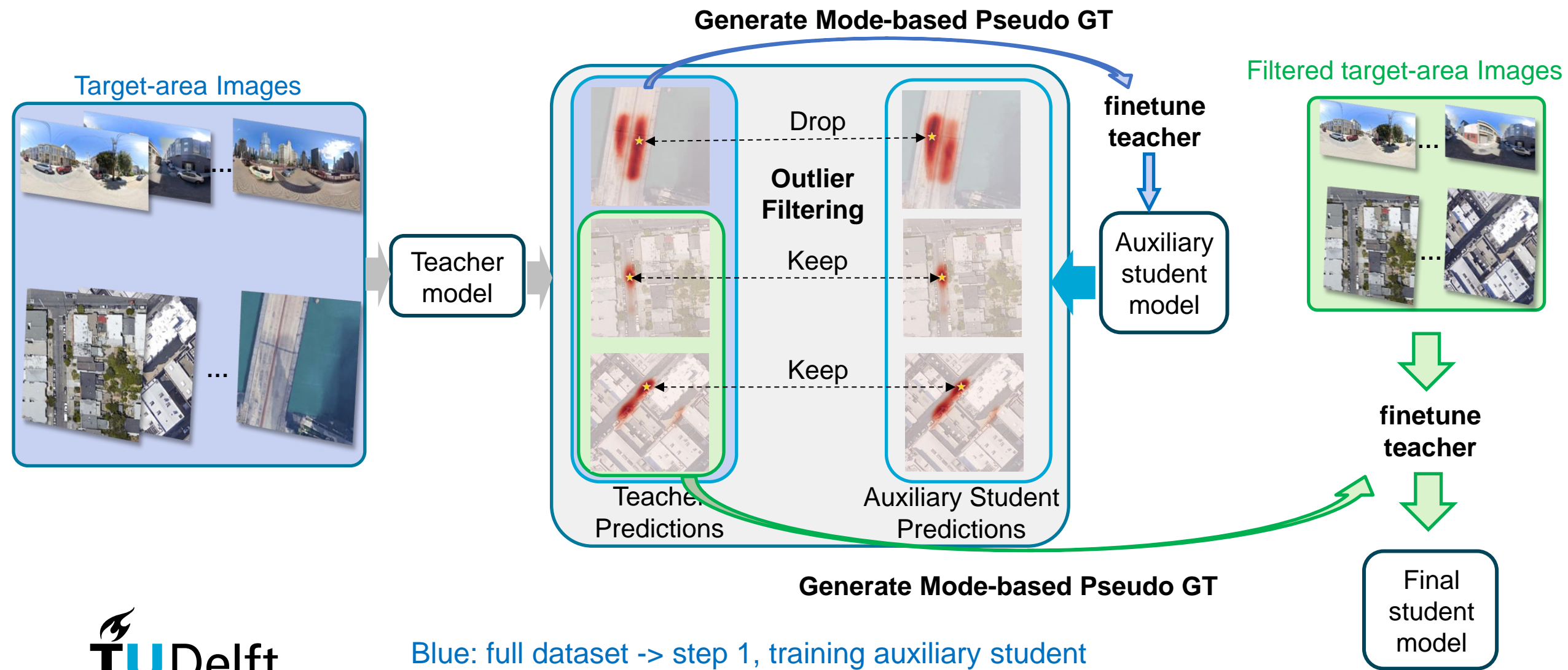
Teacher's prediction



Auxiliary student's prediction



# Leveraging noisy location data to adapt the models to new areas



# Experimental results

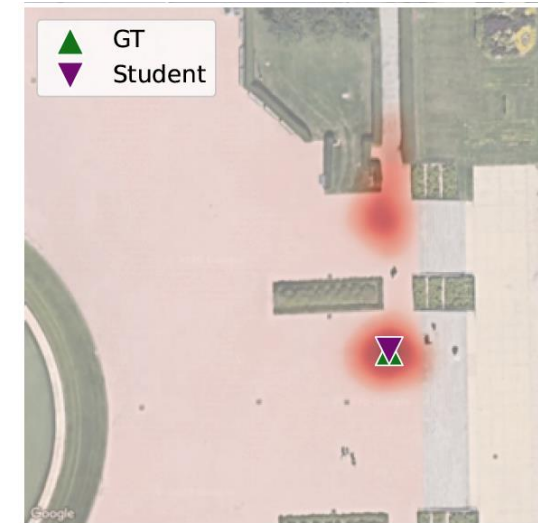
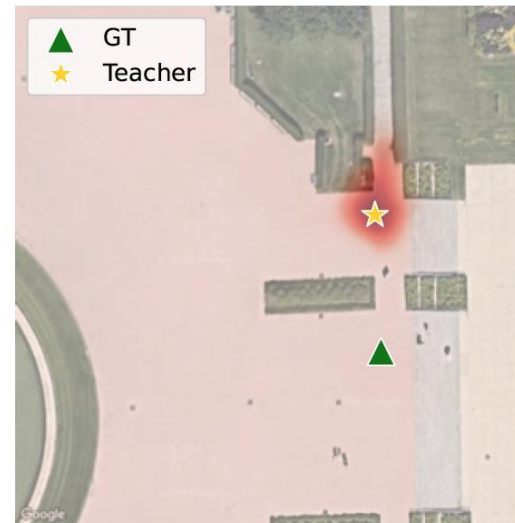
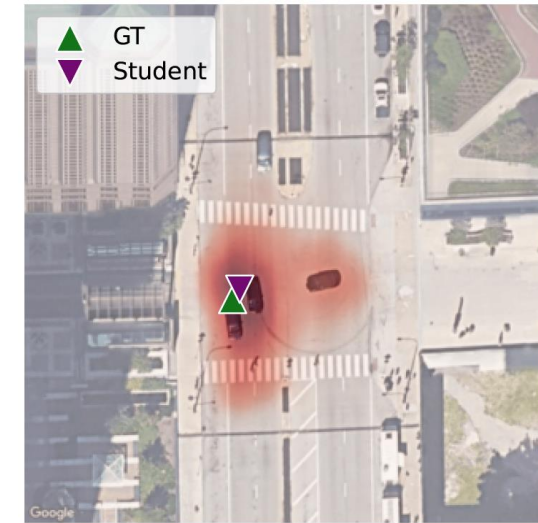
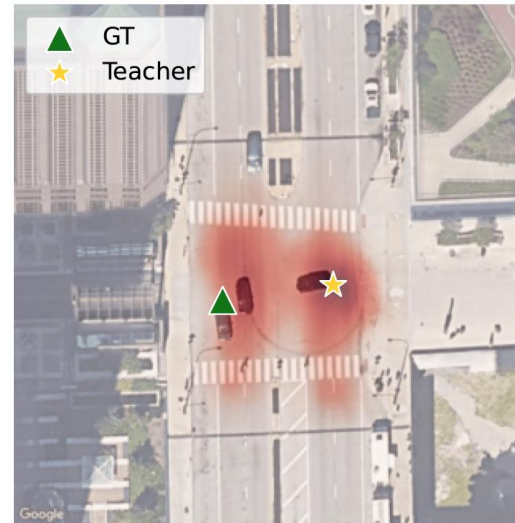
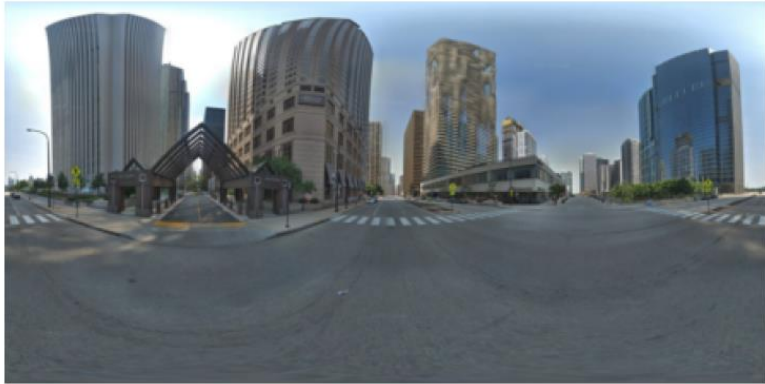
VIGOR, cross-area test	Known orientation		Unknown orientation	
	Mean (m)	Median (m)	Mean (m)	Median (m)
CCVPE [56]	4.38	1.76	5.35	1.97
CCVPE student (ours)	<b>3.85</b> (↓ 12%)	<b>1.57</b> (↓ 11%)	<b>4.27</b> (↓ 20%)	<b>1.67</b> (↓ 15%)
GGCVT [39]	5.19	1.39	-	-
GGCVT student (ours)	<b>4.34</b> (↓ 16%)	<b>1.32</b> (↓ 5%)	-	-

KITTI, cross-area test	Longitudinal error		Lateral error	
	Mean (m)	Median (m)	Mean (m)	Median (m)
CCVPE [56]	6.55	2.55	1.82	<b>0.98</b>
CCVPE student (ours)	<b>6.18</b> (↓ 6%)	<b>2.35</b> (↓ 8%)	<b>1.76</b> (↓ 3%)	<b>0.98</b> (↓ 0%)
GGCVT [39]	9.27	4.66	2.19	0.85
GGCVT student (ours)	<b>8.56</b> (↓ 8%)	<b>4.35</b> (↓ 7%)	<b>1.90</b> (↓ 13%)	<b>0.79</b> (↓ 7%)

- Our approach shows consistent and considerable improvement for two SOTA methods on two benchmarks
- We also tried generic Domain Adaptation using Entropy Minimization, but that does not work: it just makes the heat maps sharper, but does not resolve wrong modes!

# Experimental results





# Thank you! Questions?

## SliceMatch: Geometry-guided Aggregation for Cross-View Pose Estimation Lentsch et al., CVPR 2023



[Github](#)



[arXiv](#)

## Convolutional Cross-View Pose Estimation Xia et al., T-PAMI 2023



[Github](#)



[arXiv](#)

## Adapting Fine-Grained Cross-View Localization to Areas without Fine Ground Truth Xia et al., ECCV 2024



[Github](#)



[arXiv](#)

### SliceMatch: Geometry-guided Aggregation for Cross-View Pose Estimation

Ted Lentsch<sup>\*</sup> Zimin Xia<sup>\*</sup> Holger Caesar Julian F. P. Kooij  
Intelligent Vehicles Group, Delft University of Technology, The Netherlands  
{T.deVriesLentsch,Z.Xia,H.Caesar,J.F.P.Kooij}@tudelft.nl

#### Abstract

This work addresses cross-view camera pose estimation, i.e., determining the 3-Degrees-of-Freedom camera pose of a given ground-level image w.r.t. an aerial image of the local area. We propose SliceMatch, which consists of ground and aerial feature extractors, feature aggregators, and a pose predictor. The feature extractors extract dense features from the ground and aerial images. Given a set of candidate camera poses, the feature aggregators construct a single ground descriptor and a set of pose-dependent aerial descriptors. Notably, our novel aerial feature aggregator has a cross-view attention module for ground-view guided aerial feature selection and utilizes the geometric projection of the ground camera's viewing frustum on the aerial image to pool features. The efficient construction of aerial descriptors is achieved using precomputed masks. SliceMatch is trained using contrastive learning and pose estimation is formulated as a similarity comparison between the ground descriptor and the aerial descriptors. Compared to the state-of-the-art, SliceMatch achieves a 19% lower median localization error on the VIGOR benchmark using the same VGG16 backbone at 150 frames per second, and

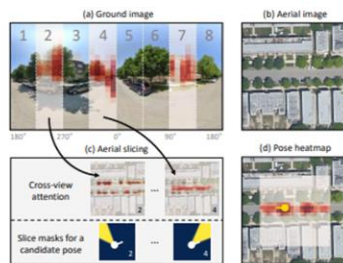


Figure 1. SliceMatch identifies for a ground-level image (a) its camera's 3-DoF pose within a corresponding aerial image (b). It divides the camera's Horizontal Field-of-View (HFOV) into 'slices', i.e., vertical regions in (a). After self-attention, our novel aggregation step (c) applies cross-view attention to create ground slice-specific aerial feature maps. To efficiently test many candidate poses, the slice features are aggregated using pose-dependent aerial slice masks that represent the camera's sliced HFOV at that pose. The slice masks for each pose are precomputed. All aerial

### Convolutional Cross-View Pose Estimation

Zimin Xia, Olaf Booi, and Julian F. P. Kooij, Member, IEEE

**Abstract**—We propose a novel end-to-end method for cross-view pose estimation. Given a ground-level query image and an aerial image that covers the query's local neighborhood, the 3 Degrees-of-Freedom camera pose of the query is estimated by matching its image descriptor to descriptors of local regions within the aerial image. The orientation-aware descriptors are obtained by using a translational equivariant convolutional ground image encoder and contrastive learning. The Localization Decoder produces a dense probability distribution in a coarse-to-fine manner with a novel Localization Matching Upsampling module. A smaller Orientation Decoder produces a vector field to condition the orientation estimate on the localization. Our method is validated on the VIGOR and KITTI datasets, where it surpasses the state-of-the-art baseline by 72% and 36% in median localization error for comparable orientation estimation accuracy. The predicted probability distribution can represent localization ambiguity, and enables rejecting possible erroneous predictions. Without re-training, the model can infer on ground images with different field of views and utilize orientation priors if available. On the Oxford RobotCar dataset, our method can reliably estimate the ego-vehicle's pose over time, achieving a median localization error under 1 meter and a median orientation error of around 1 degree at 14 FPS.

**Index Terms**—Cross-view matching, camera pose estimation, aerial imagery, localization, orientation estimation.

#### 1 INTRODUCTION

LOCALIZATION is a core task in autonomous driving and outdoor robotics [1]. In urban canyons [2], Global Navigation Satellite System (GNSS) such as GPS, often has positioning errors of up to tens of meters due to the multipath effect. Thus, other sensors [3], such as camera [4], [5] and LIDAR [6], [7], are used in combination with detailed HD maps [8], [9] to enhance the localization accuracy and robustness. In practice, most commercial vehicles are not equipped with expensive LIDAR sensors. Besides, maintaining an up-to-date HD map is laborious and expensive, especially for areas in fast development. Hence, exploring alternative map sources for camera-based methods is an important and practical task. One promising map source is aerial imagery as it provides rich semantic information with global coverage.

We consider the task of cross-view camera pose estimation, namely, estimating the camera's location and orientation from a given ground-level query image by matching it to geo-referenced aerial imagery. Previous deep learning

to have large positioning errors. A few pioneer works [25], [26], [29], [30], [31], [32], [33] demonstrated the feasibility of pinpointing the 2D location, sometimes together with the orientation, of the ground camera within a known aerial image. Similar to [28], [29], [33], we are interested in the 3-Degrees-of-Freedom (3-DoF) camera pose, i.e. planar location and orientation (yaw), instead of the full 6-DoF pose, since the change in camera height, pitch, and roll are often very small in autonomous driving.

However, several gaps must be filled before large-scale real-world deployment of cross-view camera pose estimation methods is a realistic possibility for self-driving. So far, the localization accuracy of existing methods is not yet good enough for autonomous driving requirements, e.g. the lateral and longitudinal error should be below 0.29m [34]. Besides, many methods cannot be run in real-time, i.e.  $\sim 15$  frames per second (FPS) in self-driving datasets [35], [37], because of using expensive iterative optimization [28], [30] or computationally heavy Transformers [33]. We also

### Adapting Fine-Grained Cross-View Localization to Areas without Fine Ground Truth

Zimin Xia<sup>1</sup>, Yujiao Shi<sup>2</sup>, Hongdong Li<sup>3</sup>, and Julian F. P. Kooij<sup>4</sup>

<sup>1</sup> École Polytechnique Fédérale de Lausanne (EPFL), Switzerland

zimin.xia@epfl.ch

<sup>2</sup> ShanghaiTech University, China

<sup>3</sup> Australian National University, Australia

<sup>4</sup> Delft University of Technology, The Netherlands

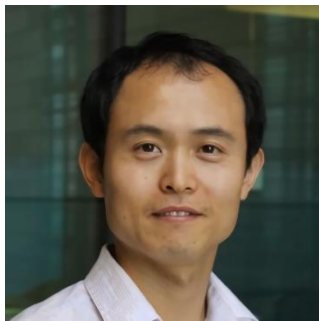
**Abstract.** Given a ground-level query image and a geo-referenced aerial image that covers the query's local surroundings, fine-grained cross-view localization aims to estimate the location of the ground camera inside the aerial image. Recent works have focused on developing advanced networks trained with accurate ground truth (GT) locations of ground images. However, the trained models always suffer a performance drop when applied to images in a new target area that differs from training. In most deployment scenarios, acquiring fine GT, i.e. accurate GT locations, for target-area images to re-train the network can be expensive and sometimes infeasible. In contrast, collecting images with noisy GT with errors of tens of meters is often easy. Motivated by this, our paper focuses on improving the performance of a trained model in a new target area by leveraging only the target-area images without fine GT. We propose a weakly supervised learning approach based on knowledge self-distillation. This approach uses predictions from a pre-trained model as pseudo GT to supervise a copy of itself. Our approach includes a mode-based pseudo GT generation for reducing uncertainty in pseudo GT and an outlier filtering method to remove unreliable pseudo GT. Our approach is validated using two recent state-of-the-art models on two benchmarks. The results demonstrate that it consistently and considerably boosts the localization accuracy in the target area.

# On the Estimation of Image-matching Uncertainty in Visual Place Recognition

CVPR 2024 (*Poster Highlight*)



Mubariz Zaffar



Liangliang Nan



Julian F. P. Kooij



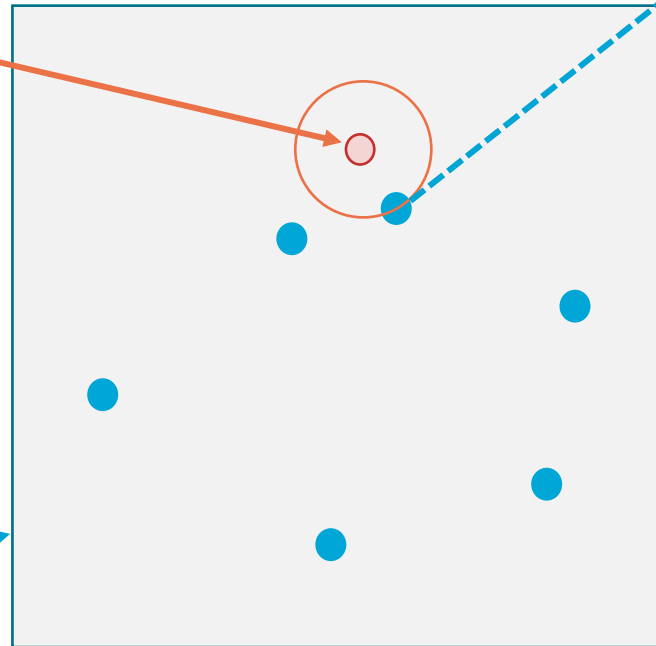
# Visual Place Recognition (VPR)

Where is this query?



CNN

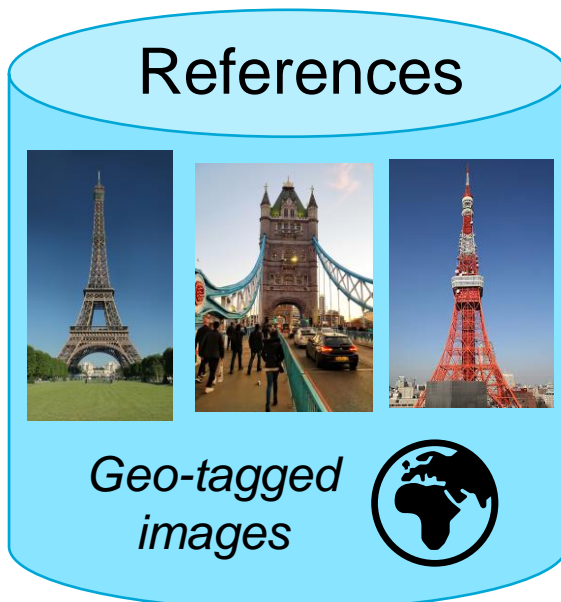
Feature space



CNN



Retrieve  
Nearest neighbour

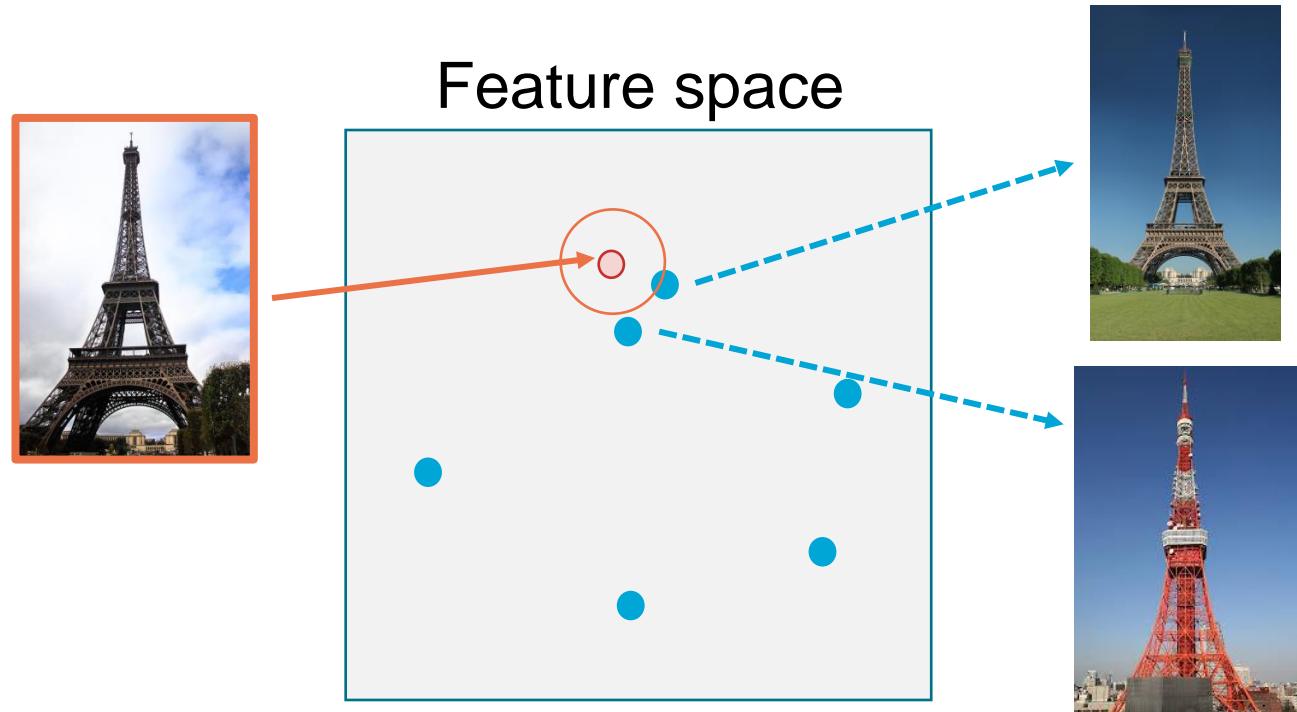


**True Positive:** NN is located within  $x$  meters from query's true location



# VPR Matching Uncertainty

- How do we know if retrieved VPR result is reliable?
- “Visual aliasing”: some locations just look similar (patch of sky, white wall, ...)



- Various procedures to estimate VPR Matching Uncertainty exist ...

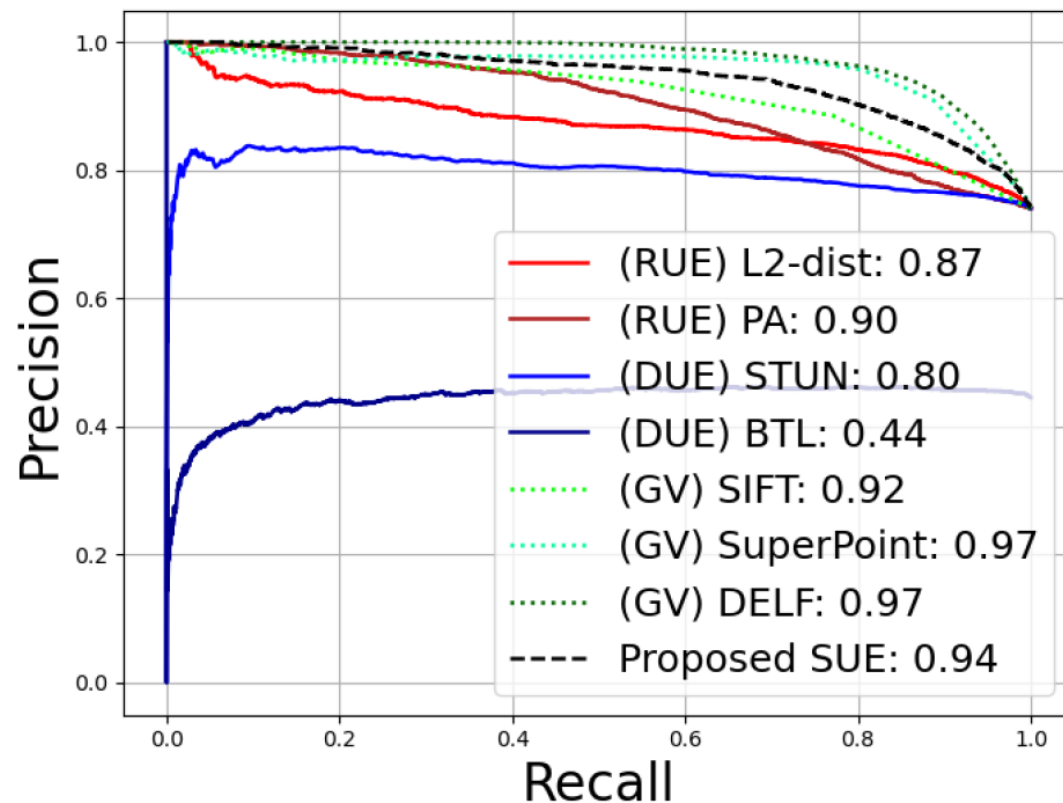
# VPR Matching Uncertainty

We categorize current VPR uncertainty estimation

- **RUE**: Retrieval-based uncertainty estimation
- **DUE**: Data-driven uncertainty estimation
- **GV**: Geometric verification

Propose a new simple approach

- **SUE**: Spatial Uncertainty Estimation

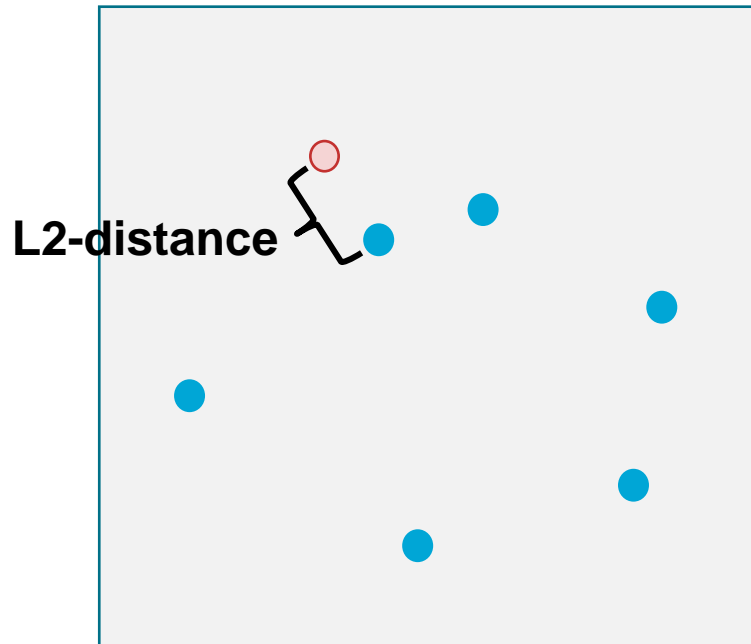


# RUE: Retrieval-based uncertainty estimation

Common confidence score used in Visual SLAM for loop closure

- **L2-distance:** between query feature and best-matching feature
- **PA-Score:** ratio of L2-distance between 1<sup>st</sup> and 2<sup>nd</sup> nearest neighbour reference

*Feature space*



Least uncertain



Most uncertain





# DUE: Data-driven uncertainty estimation

Estimate (aleatoric) uncertainty from image content

- **Bayesian Triplet Loss (BTL):** Warburg et al., ICCV 2021
- **STUN:** Cai et al., IROS 2022

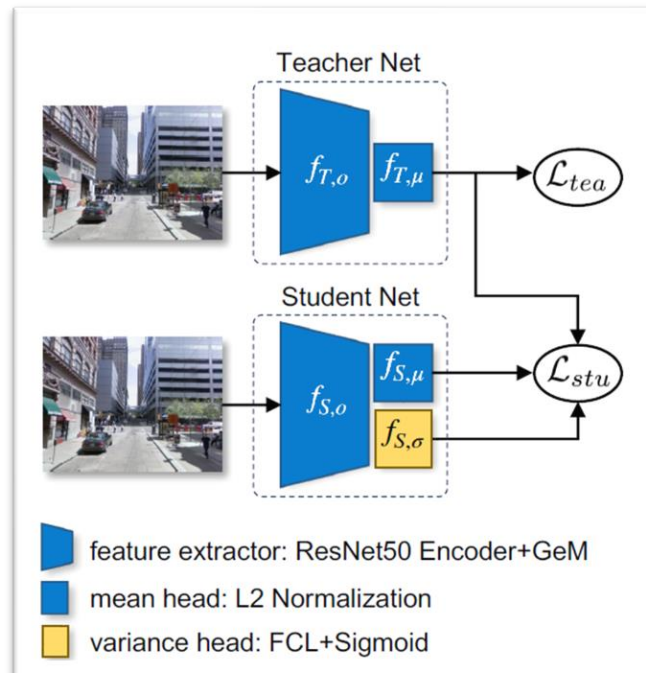
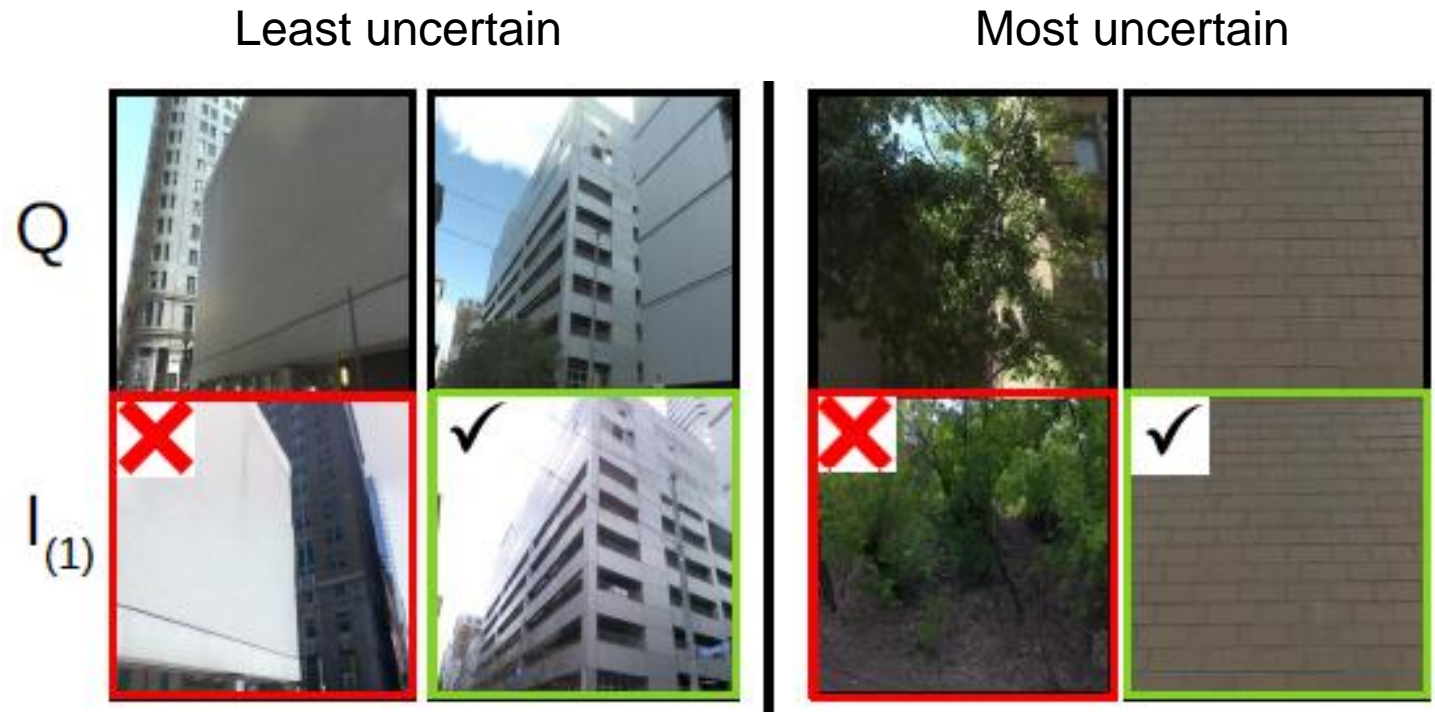


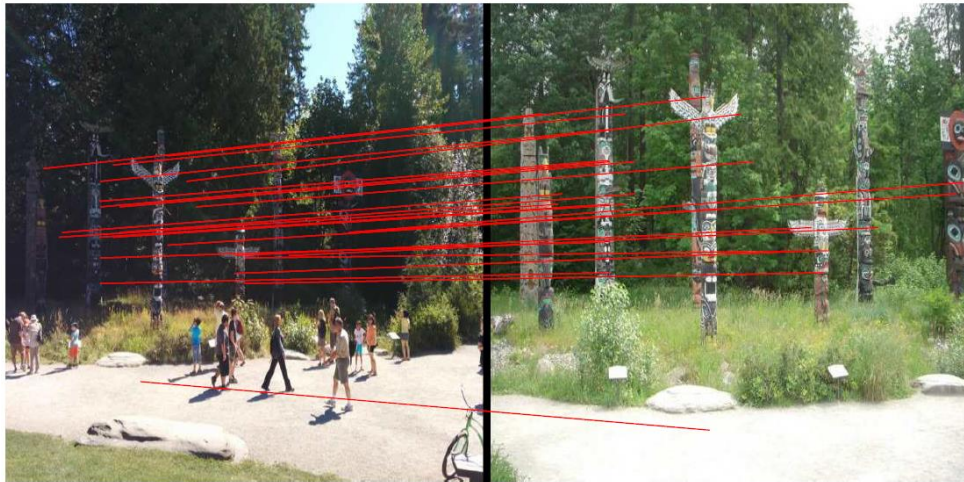
Image: [Cai et al., IROS 2022]



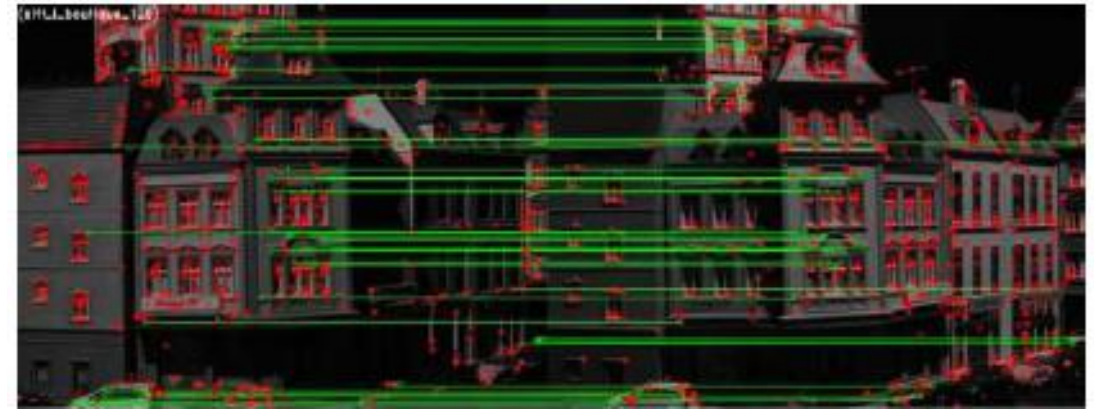
# GV: Geometric Verification

Local feature matching + RANSAC  
Slow but accurate, typically used to rerank top-K retrieved candidates

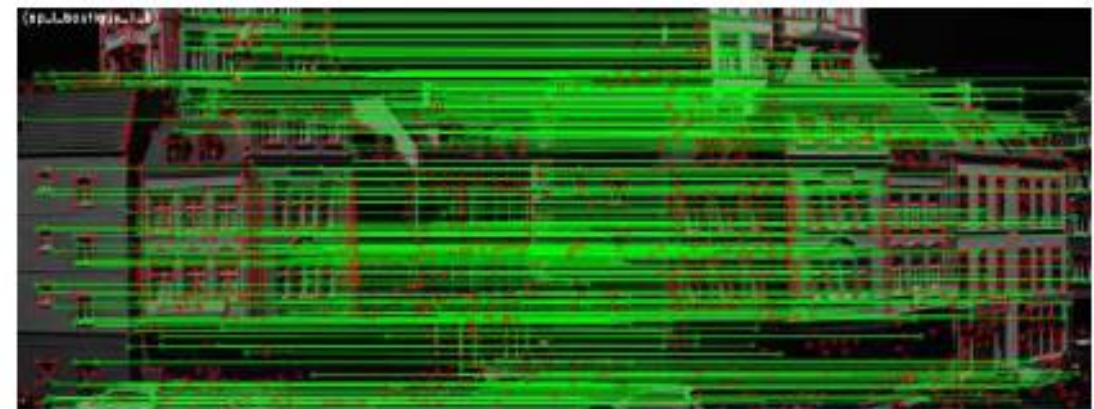
- **SIFT**: Lowe, IJCV 2004
- **DELf**: Noh et al., ICCV 2017
- **SuperPoint**: DeTone, CVPRW 2018



DELf (image credit: Noh, ICCV'17)



SIFT (image credit: [DeTone, CVPRW'18])



SuperPoint (image credit: DeTone, CVPRW'18)

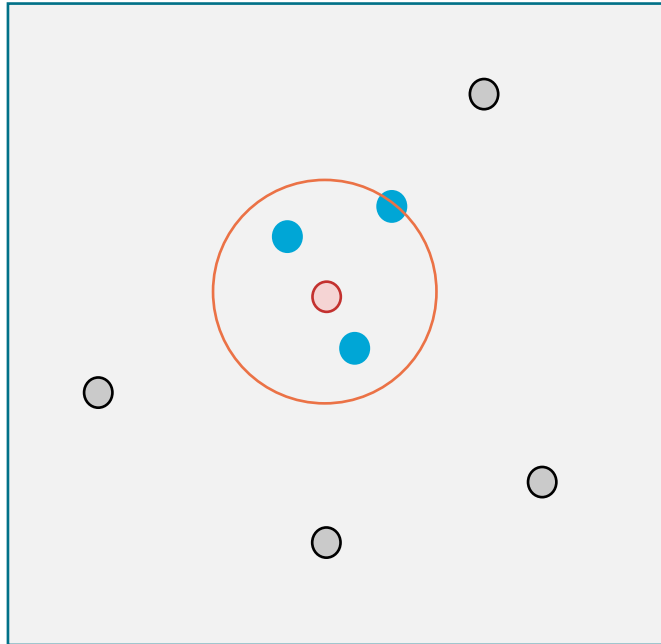
# SUE: Spatial Uncertainty Estimation

- Existing approaches do not consider geo-locations of the references!
- SUE estimates if visual aliasing occurs between distant map locations
- *SUE is ridiculously simple!*
  1. For a query, retrieve top-k best matching references
  2. Weigh top-k references inversely by L2-distance
  3. Compute weighted variance of *spatial locations* of top-k references
  4. Total variance is matching uncertainty score



# SUE: Spatial Uncertainty Estimation

*Feature space*



Spatial spread  
of top-K references

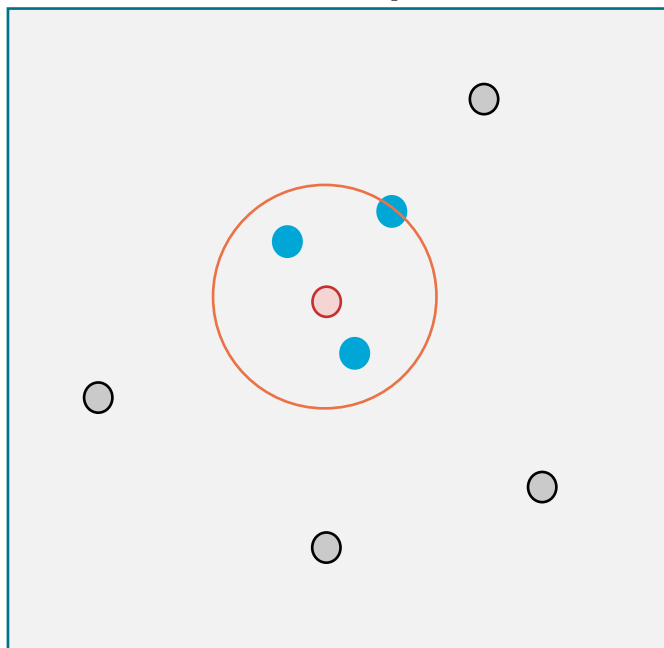


Query

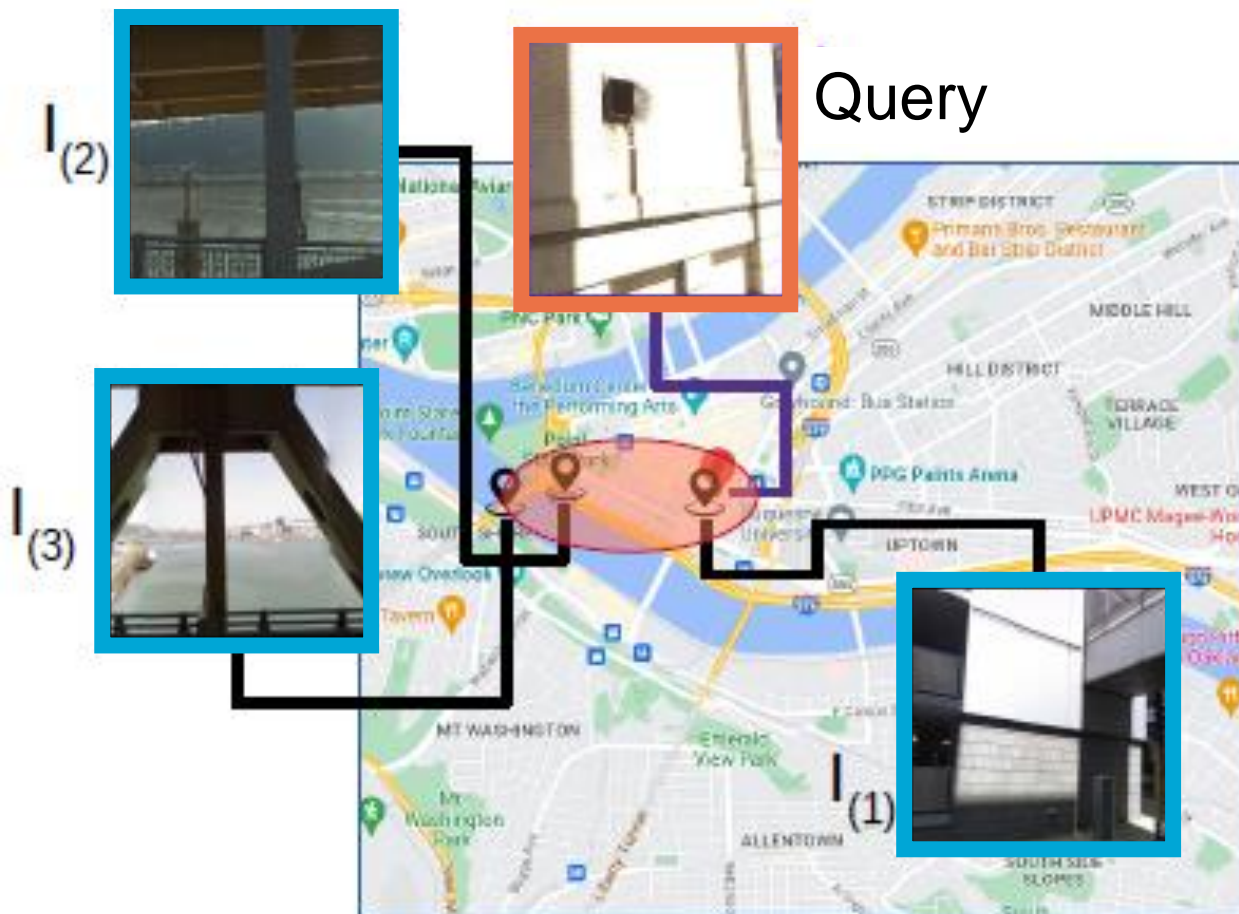
Low spatial spread  $\rightarrow$  low SUE uncertainty

# SUE: Spatial Uncertainty Estimation

*Feature space*

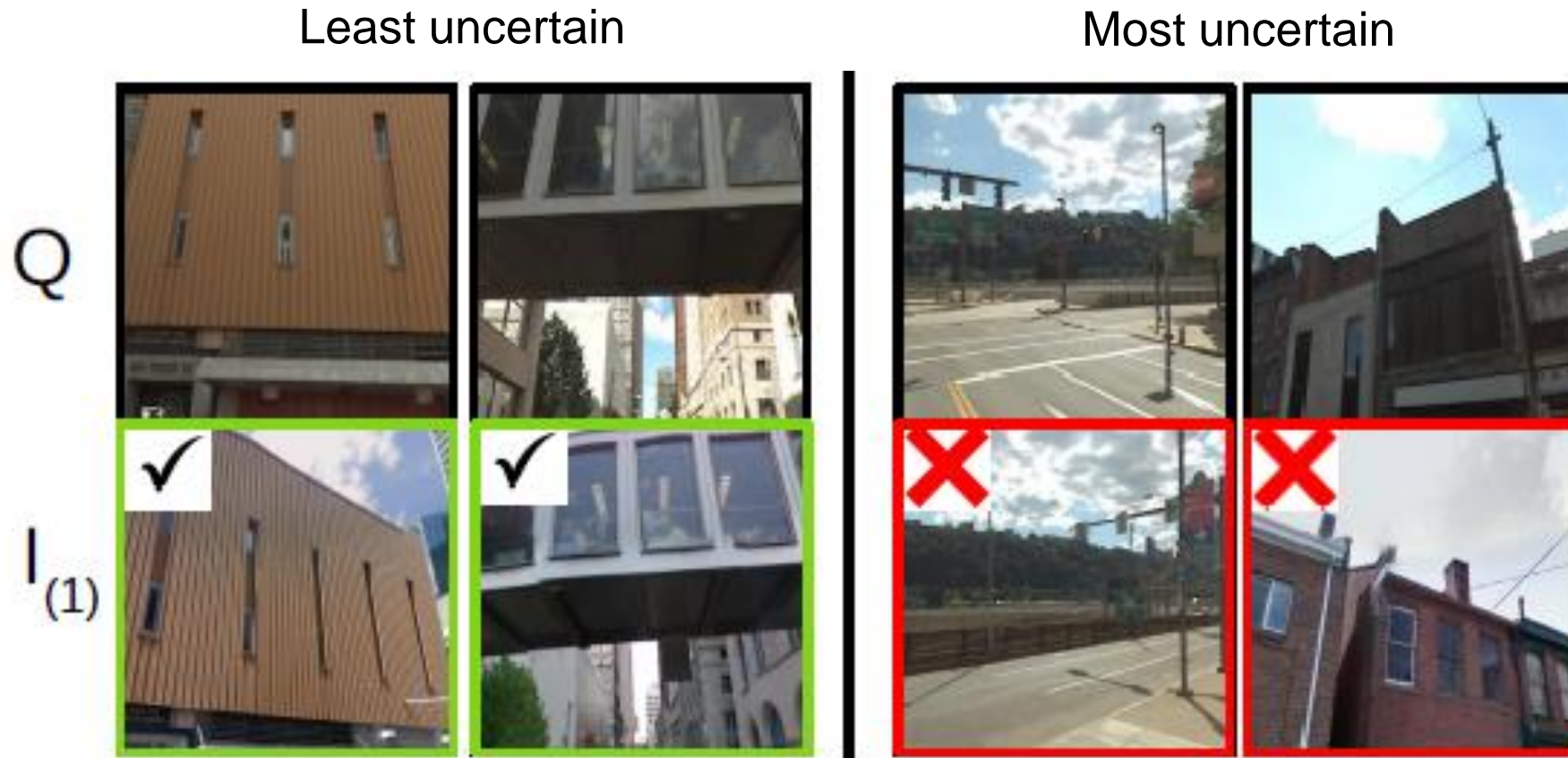


Spatial spread  
of top-K references



High spatial spread  $\rightarrow$  high SUE uncertainty

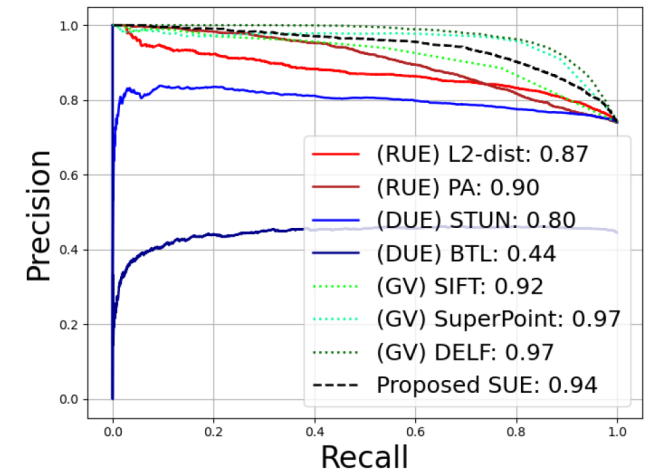
# SUE: Spatial Uncertainty Estimation





# Results

- AUC-PR on various VPR benchmarks
- Processing time (milliseconds)



Method	↑ Pitts.	↑ Sanfr.	↑ Stluc.	↑ Eyn.	↑ MSLS	↑ Nordland	↑ Average	↓ Time
(RUE) L2-distance	0.87	0.76	0.79	0.87	0.64	0.18	0.69	<b>0.05</b>
(RUE) PA-Score [18]	0.90	0.65	0.77	0.88	0.68	0.21	0.68	<b>0.05</b>
(DUE) BTL [50]	0.44	0.17	0.34	0.45	0.21	0.07	0.28	0.20
(DUE) STUN [9]	0.79	0.57	0.66	0.71	0.44	0.05	0.54	0.10
SUE	<b>0.94</b>	<b>0.84</b>	<b>0.88</b>	<b>0.93</b>	<b>0.77</b>	<b>0.26</b>	<b>0.77</b>	1.08
(GV) SIFT-RANSAC [27]	0.92	0.89	0.93	<b>0.96</b>	0.70	0.15	0.76	<b>129</b>
(GV) DELF-RANSAC [33]	0.97	<b>0.92</b>	<b>0.97</b>	0.95	<b>0.95</b>	<b>0.84</b>	<b>0.93</b>	1587
(GV) Super-RANSAC [15]	0.95	<b>0.95</b>	<b>0.97</b>	<b>0.96</b>	0.87	0.50	0.87	848

# Examples



Same location  
Query matched  
Low uncertainty  
desirable

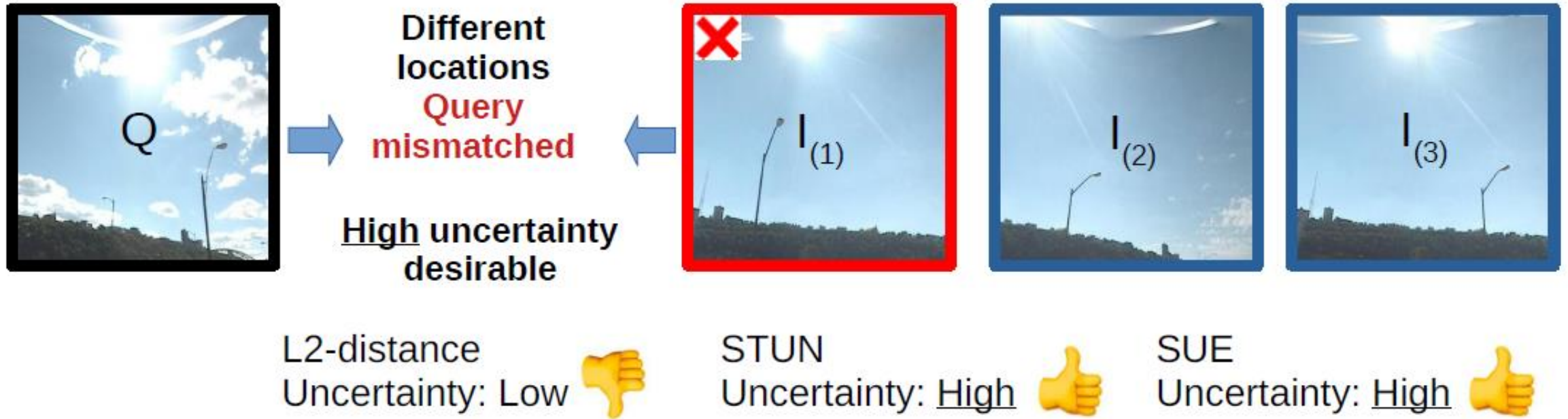


L2-distance  
Uncertainty: Low 👍

STUN  
Uncertainty: High 👎

SUE  
Uncertainty: Low 👍

# Examples





# Conclusions

- Compared different approaches for estimating VPR uncertainty
- Recent work tends to ignore simple baselines!
- Simple L2-distance already outperforms Deep Learning-based
- We propose SUE, which also looks at reference locations
- SUE is best efficient method
- In paper we show SUE can complement computationally expensive Geometric Verification

# Thank you! Questions?

## *SliceMatch: Geometry-guided Aggregation for Cross-View Pose Estimation* Lentsch et al., CVPR 2023



[Github](#)



[arXiv](#)

## *Convolutional Cross-View Pose Estimation* Xia et al., T-PAMI 2023



[Github](#)



[arXiv](#)

## *On the Estimation of Image-matching Uncertainty in Visual Place Recognition* Zaffar et al., CVPR 2024



[arXiv](#)

### SliceMatch: Geometry-guided Aggregation for Cross-View Pose Estimation

Ted Lentsch\* Zimin Xia\* Holger Caesar Julian F. P. Kooij  
Intelligent Vehicles Group, Delft University of Technology, The Netherlands  
{T.deVriesLentsch,Z.Xia,H.Caesar,J.F.P.Kooij}@tudelft.nl

#### Abstract

This work addresses cross-view camera pose estimation, i.e., determining the 3-Degrees-of-Freedom camera pose of a given ground-level image w.r.t. an aerial image of the local area. We propose SliceMatch, which consists of ground and aerial feature extractors, feature aggregators, and a pose predictor. The feature extractors extract dense features from the ground and aerial images. Given a set of candidate camera poses, the feature aggregators construct a single ground descriptor and a set of pose-dependent aerial descriptors. Notably, our novel aerial feature aggregator has a cross-view attention module for ground-view guided aerial feature selection and utilizes the geometric projection of the ground camera's viewing frustum on the aerial image to pool features. The efficient construction of aerial descriptors is achieved using precomputed masks. SliceMatch is trained using contrastive learning and pose estimation is formulated as a similarity comparison between the ground descriptor and the aerial descriptors. Compared to the state-of-the-art, SliceMatch achieves a 19% lower median localization error on the VIGOR benchmark using the same VGG16 backbone at 150 frames per second, and

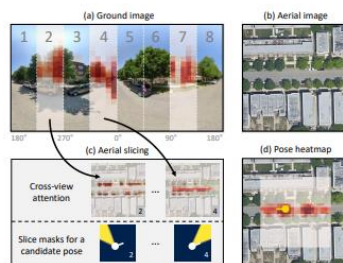


Figure 1. SliceMatch identifies for a ground-level image (a) its camera's 3-DoF pose within a corresponding aerial image (b). It divides the camera's Horizontal Field-of-View (HFOV) into 'slices', i.e., vertical regions in (a). After self-attention, our novel aggregation step (c) applies cross-view attention to create ground slice-specific aerial feature maps. To efficiently test many candidate poses, the slice features are aggregated using pose-dependent aerial slice masks that represent the camera's sliced HFOV at that pose. The slice masks for each pose are precomputed. All aerial

### Convolutional Cross-View Pose Estimation

Zimin Xia, Olaf Booij, and Julian F. P. Kooij, Member, IEEE

**Abstract**—We propose a novel end-to-end method for cross-view pose estimation. Given a ground-level query image and an aerial image that covers the query's local neighborhood, the 3 Degrees-of-Freedom camera pose of the query is estimated by matching its image descriptor to descriptors of local regions within the aerial image. The orientation-aware descriptors are obtained by using a translational equivariant convolutional ground image encoder and contrastive learning. The Localization Decoder produces a dense probability distribution in a coarse-to-fine manner with a novel Localization Matching Upsampling module. A smaller Orientation Decoder produces a vector field to condition the orientation estimate on the localization. Our method is validated on the VIGOR and KITTI datasets, where it surpasses the state-of-the-art baseline by 72% and 36% in median localization error for comparable orientation estimation accuracy. The predicted probability distribution can represent localization ambiguity, and enables rejecting possible erroneous predictions. Without re-training, the model can infer on ground images with different field of views and utilize orientation priors if available. On the Oxford RobotCar dataset, our method can reliably estimate the ego-vehicle's pose over time, achieving a median localization error under 1 meter and a median orientation error of around 1 degree at 14 FPS.

**Index Terms**—Cross-view matching, camera pose estimation, aerial imagery, localization, orientation estimation.

#### 1 INTRODUCTION

LOCALIZATION is a core task in autonomous driving and outdoor robotics [1]. In urban canyons [2], Global Navigation Satellite System (GNSS) such as GPS, often has positioning errors of up to tens of meters due to the multipath effect. Thus, other sensors [3], such as camera [4], [5] and LIDAR [6], [7], are used in combination with detailed HD maps [8], [9] to enhance the localization accuracy and robustness. In practice, most commercial vehicles are not equipped with expensive LIDAR sensors. Besides, maintaining an up-to-date HD map is laborious and expensive, especially for areas in fast development. Hence, exploring alternative map sources for camera-based methods is an important and practical task. One promising map source is aerial imagery as it provides rich semantic information with global coverage.

We consider the task of cross-view camera pose estimation, namely, estimating the camera's location and orientation from a given ground-level query image by matching it to geo-referenced aerial imagery. Previous deep learning

to have large positioning errors. A few pioneer works [2], [29], [30], [31], [32], [33] demonstrated the feasibility of pinpointing the 2D location, sometimes together with the orientation, of the ground camera within a known aerial image. Similar to [29], [33], we are interested in the 3-Degrees-of-Freedom (3-DoF) camera pose, i.e. planar location and orientation (yaw), instead of the full 6-DoF pose, since the change in camera height, pitch, and roll are often very small in autonomous driving.

However, several gaps must be filled before large-scale real-world deployment of cross-view camera pose estimation methods is a realistic possibility for self-driving. So far, the localization accuracy of existing methods is not yet good enough for autonomous driving requirements, e.g. the lateral and longitudinal error should be below 0.29m [34]. Besides, many methods cannot be run in real-time, i.e.  $\sim 15$  frames per second (FPS) in self-driving datasets [35], [36], because of using expensive iterative optimization [29], [37] or computationally heavy Transformers [33]. We also

### On the Estimation of Image-matching Uncertainty in Visual Place Recognition

Mubariz Zaffar  
ME, TU Delft  
The Netherlands

M.Zaffar@tudelft.nl

Liangliang Nan  
ABE, TU Delft  
The Netherlands

Liangliang.Nan@tudelft.nl

Julian F. P. Kooij  
ME, TU Delft  
The Netherlands

J.F.P.Kooij@tudelft.nl

#### Abstract

In Visual Place Recognition (VPR) the pose of a query image is estimated by comparing the image to a map of reference images with known reference poses. As is typical for image retrieval problems, a feature extractor maps the query and reference images to a feature space, where a nearest neighbor search is then performed. However, till recently little attention has been given to quantifying the confidence that a retrieved reference image is a correct match. Highly certain but incorrect retrieval can lead to catastrophic failure of VPR-based localization pipelines. This work compares for the first time the main approaches for estimating the image-matching uncertainty, including the traditional retrieval-based uncertainty estimation, more recent data-driven aleatoric uncertainty estimation, and the compute-intensive geometric verification. We further formulate a simple baseline method, "SUE", which unlike the other methods considers the freely-available poses of the reference images in the map. Our experiments reveal that

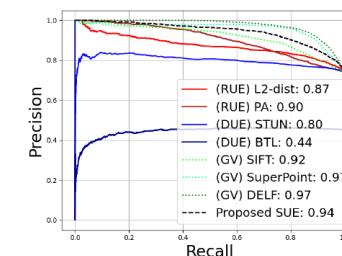


Figure 1. The Precision-Recall curves on the Pittsburgh dataset [4] for the three common categories of VPR uncertainty estimation methods (RUE, DUE, GV), and for our proposed baseline SUE which uniquely considers spatial locations of the top-K references. The global image descriptors [9] are fixed for all methods except