# Weakly-supervised Camera Localization by Ground-to-satellite Image Registration

Yujiao Shi<sup>1</sup><sup>(6)</sup>, Hongdong Li<sup>2</sup><sup>(6)</sup>, Akhil Perincherry<sup>3</sup><sup>(6)</sup>, and Ankit Vora<sup>3</sup><sup>(6)</sup>

<sup>1</sup> ShanghaiTech University, China
 <sup>2</sup> The Australian National University, Australia
 <sup>3</sup> Ford Motor Company, USA
 shiyj2@shanghaitech.edu.cn

## 1 Sensitiveness to Coarse Pose Errors

In this section, we investigate the performance of our method under varying coarse pose errors.

**Range of location errors.** Table A presents the performance comparison between our method and the state-of-the-art [26, 28], across different ranges of location errors:  $28 \times 28 \text{ m}^2$  and  $56 \times 56 \text{ m}^2$ , while maintaining the same orientation ambiguity of 20°. The results show that our method achieves consistently the best performance on cross-area evaluation.

$\begin{array}{c} Location \\ Error (m^2) \end{array}$	Algorithms	$\begin{vmatrix} & \text{La} \\ d = 1 \uparrow \end{vmatrix}$	$\begin{array}{c} \text{teral} \\ d = 3 \uparrow \end{array}$	Test-1 (S Longi $d = 1 \uparrow$	$\begin{array}{l} \text{fame-area} \\ \text{tudinal} \\ d = 3 \uparrow \end{array}$	$\begin{array}{c} \text{Azimuth} \\ \theta = 1 \uparrow  \theta = 3 \uparrow \end{array}$		
$28 \times 28$	Shi and Li $[26]^*$ Shi <i>et al.</i> $[28]^*$	44.66 85.85	$73.92 \\ 98.46$	12.06 23.27	$35.62 \\ 46.99$	25.31 98.89	$57.41 \\ 99.97$	
20 / 20	Ours $(\lambda = 0)$ Ours $(\lambda = 1)$	61.46 66.39	$87.76 \\ 94.38$	13.44 18.18	$38.14 \\ 53.59$	99.76 99.76	$\begin{array}{c} 100.00\\ 100.00 \end{array}$	
$56 \times 56$	Shi and Li $[26]^*$ Shi <i>et al.</i> $[28]^*$	$\begin{vmatrix} 35.54 \\ 76.44 \end{vmatrix}$	70.77 96.34	5.22 23.54	$15.88 \\ 50.57$	19.64 99.10	$51.76 \\ 100.00$	
	$\begin{vmatrix} \mathbf{Ours} & (\lambda = 0) \\ \mathbf{Ours} & (\lambda = 1) \end{vmatrix}$	59.58 66.07	$85.74 \\ 94.22$	$  11.37 \\ 16.51$	$\begin{array}{c} 31.94 \\ 49.96 \end{array}$	$99.66 \\ 99.66$	$\begin{array}{c} 100.00\\ 100.00\end{array}$	
$\begin{array}{c} \text{Location} \\ \text{Error (m}^2) \end{array}$	Algorithms	$\begin{vmatrix} & \text{La} \\ d = 1 \uparrow \end{vmatrix}$	$\begin{array}{c} \text{teral} \\ d = 3 \uparrow \end{array}$	$\begin{array}{c c} \textbf{Test-2} & \textbf{(C)} \\ & \text{Longi} \\ d = 1 \uparrow \end{array}$	Cross-area) tudinal $d = 3 \uparrow$	$\begin{vmatrix} & \text{Azin} \\ \theta = 1 \uparrow \end{vmatrix}$	$\begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \end{array}$	
Location Error $(m^2)$ $28 \times 28$	Algorithms Shi and Li [26]* Shi et al. [28]*	$\begin{vmatrix} & \text{La} \\ d = 1 \uparrow \\ & 34.17 \\ 60.01 \end{vmatrix}$	teral $d = 3 \uparrow$ 72.30 87.96	Test-2 (C         Longi $d = 1 \uparrow$ 11.56         14.69	$\begin{array}{c} \textbf{Cross-area}\\ \textbf{tudinal}\\ d=3\uparrow\\ \hline\\ 35.08\\ 35.64 \end{array}$	$\begin{vmatrix} Azin\\ \theta = 1 \uparrow \\ 11.40\\ 99.42 \end{vmatrix}$	$ \begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \\ \hline 48.18 \\ 100.00 \end{array} $	
$\frac{\text{Location}}{\text{Error (m}^2)}$ $28 \times 28$	AlgorithmsShi and Li [26]*Shi et al. [28]*Ours ( $\lambda = 0$ )Ours ( $\lambda = 1$ )	$\begin{vmatrix} & \text{La:} \\ d = 1 \uparrow \\ & 34.17 \\ 60.01 \\ & 65.62 \\ & 67.90 \end{vmatrix}$	teral $d = 3 \uparrow$ 72.30 87.96 90.32 89.76	$\begin{array}{c c} \textbf{Test-2} (\textbf{C} \\ Longi \\ d = 1 \uparrow \\ 11.56 \\ 14.69 \\ 13.46 \\ 14.29 \end{array}$	<b>Cross-area</b> tudinal $d = 3 \uparrow$ 35.08         35.64         38.53         42.92	$\begin{vmatrix} Azin\\ \theta = 1 \uparrow\\ 11.40\\ 99.42\\ \end{vmatrix}$ $\begin{vmatrix} 99.97\\ 99.97\\ \end{vmatrix}$	muth $\theta = 3 \uparrow$ 48.18 100.00 100.00 100.00	
$ \begin{array}{c} \text{Location} \\ \text{Error } (m^2) \\ \hline 28 \times 28 \\ \hline 56 \times 56 \\ \end{array} $	AlgorithmsShi and Li [26]*Shi et al. [28]*Ours ( $\lambda = 0$ )Ours ( $\lambda = 1$ )Shi and Li [26]*Shi et al. [28]*	$\begin{vmatrix} & \text{La:} \\ d = 1 \uparrow \\ & 34.17 \\ 60.01 \\ & 65.62 \\ 67.90 \\ & 27.82 \\ 57.72 \end{vmatrix}$	teral $d = 3 \uparrow$ 72.30 87.96 90.32 89.76 59.79 86.77	Test-2 (C           Longi $d = 1 \uparrow$ 11.56           14.69           13.46           14.29           5.75           14.15	Pross-area         tudinal $d = 3 \uparrow$ 35.08         35.64         38.53         42.92         16.36         34.59	$\begin{vmatrix} Azin\\ \theta = 1 \uparrow\\ 11.40\\ 99.42\\ 99.97\\ 99.97\\ 18.42\\ 98.98\\ \end{vmatrix}$	$ \begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \\ \hline \\ 48.18 \\ 100.00 \\ \hline \\ 100.00 \\ \hline \\ 49.72 \\ 100.00 \\ \end{array} $	

 
 Table A: Performance comparison with different location error ranges on the crossview KITTI dataset.

Orientation Ambiguity	Algorithms	$\begin{vmatrix} & \text{La} \\ d = 1 \uparrow \end{vmatrix}$	teral $d = 3 \uparrow$	Test-1 (S Longi $d = 1 \uparrow$	dame-area tudinal $d = 3 \uparrow$	$\begin{vmatrix} & \text{Azi} \\ \theta = 1 \uparrow \end{vmatrix}$	$\begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \end{array}$
20°	Shi and Li [26]* Shi <i>et al.</i> [28]*	35.54 76.44	70.77 96.34	5.22 23.54	$15.88 \\ 50.57$	19.64 99.10	$51.76 \\ 100.00$
	Ours $(\lambda = 0)$ Ours $(\lambda = 1)$	59.58 66.07	$85.74 \\ 94.22$	$ \begin{array}{c c} 11.37 \\ 16.51 \end{array} $	$\begin{array}{c} 31.94 \\ 49.96 \end{array}$	99.66 99.66	$100.00 \\ 100.00$
80°	Shi and Li [26]* Shi <i>et al.</i> [28]*	26.95 70.21	$62.39 \\ 95.47$	5.14 22.29	$\begin{array}{c} 15.69 \\ 48.90 \end{array}$	3.10 53.27	$8.88 \\93.98$
80	Ours $(\lambda = 0)$ Ours $(\lambda = 1)$	53.11 57.94	$86.03 \\ 91.49$	12.99 17.73	$32.18 \\ 47.44$	57.65 57.65	$96.79 \\ 96.79$
Orientation Ambiguity	Algorithms	$\begin{array}{c} \text{Lat} \\ d = 1 \uparrow \end{array}$	teral $d = 3 \uparrow$	Test-2 (C Longit $d = 1 \uparrow$	$\begin{array}{c} \textbf{ross-area} \\ \textbf{tudinal} \\ d = 3 \uparrow \end{array}$	$\begin{vmatrix} & \text{Azin} \\ \theta = 1 \uparrow \end{vmatrix}$	$\begin{array}{l} \text{muth} \\ \theta = 3 \uparrow \end{array}$
Orientation Ambiguity 20°	Algorithms Shi and Li [26]* Shi et al. [28]*	$\begin{array}{c} \text{Lat} \\ d = 1 \uparrow \\ 27.82 \\ 57.72 \end{array}$	teral $d = 3 \uparrow$ 59.79 86.77	<b>Test-2 (C</b> Longit $d = 1 \uparrow$ 5.75 14.15	$\begin{array}{c} \textbf{ross-area} \\ \textbf{tudinal} \\ d = 3 \uparrow \\ \hline 16.36 \\ 34.59 \end{array}$	$\begin{vmatrix} Azin\\ \theta = 1 \uparrow \\ 18.42\\ 98.98 \end{vmatrix}$	$ \begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \\ \hline 49.72 \\ 100.00 \end{array} $
Orientation Ambiguity 20°	AlgorithmsShi and Li [26]*Shi et al. [28]*Ours ( $\lambda = 0$ )Ours ( $\lambda = 1$ )	$\begin{array}{c c} & \text{Lat} \\ d = 1 \uparrow \\ 27.82 \\ 57.72 \\ 62.73 \\ 64.74 \end{array}$	teral $d = 3 \uparrow$ 59.79 86.77 86.53 86.18	Test-2 (C           Longit $d = 1 \uparrow$ 5.75           14.15           9.98           11.81	ross-area         tudinal $d = 3 \uparrow$ 16.36         34.59         29.67         34.77	$\begin{vmatrix} & \text{Azin} \\ \theta = 1 \uparrow \\ & 18.42 \\ & 98.98 \\ \end{vmatrix}$	$ \begin{array}{c} \text{muth} \\ \theta = 3 \uparrow \\ \hline 49.72 \\ 100.00 \\ \hline 100.00 \\ 100.00 \\ \end{array} $
Orientation Ambiguity 20° 80°	AlgorithmsShi and Li [26]*Shi et al. [28]*Ours ( $\lambda = 0$ )Ours ( $\lambda = 1$ )Shi and Li [26]*Shi et al. [28]*	$\begin{array}{c c} & \text{Lat} \\ d = 1 \uparrow \\ 27.82 \\ 57.72 \\ 62.73 \\ 64.74 \\ 22.43 \\ 56.97 \end{array}$	teral $d = 3 \uparrow$ 59.79 86.77 86.53 86.18 54.63 87.72	Test-2 (C           Longit $d = 1 \uparrow$ 5.75           14.15           9.98           11.81           5.17           15.17	ross-area) tudinal $d = 3 \uparrow$ 16.36 34.59 29.67 34.77 15.78 35.39	$\begin{vmatrix} Azin \\ \theta = 1 \uparrow \\ 18.42 \\ 98.98 \\ 99.99 \\ 99.99 \\ 99.99 \\ 3.05 \\ 58.68 \end{vmatrix}$	$\begin{aligned} & \text{muth} \\ \theta &= 3 \uparrow \\ & 49.72 \\ 100.00 \\ & 100.00 \\ & 100.00 \\ & 8.50 \\ & 95.92 \end{aligned}$

 
 Table B: Performance comparison with different location error ranges on the crossview KITTI dataset.

Variation in orientation ambiguity. Subsequently, we augment the orientation ambiguity from 20° to 80°, while maintaining a location error range of  $56 \times 56$  m<sup>2</sup>. Table B provides the performance comparison between our method and the two state-of-the-art [26, 28]. Our method achieves the best cross-area evaluation performance on the different orientation ambiguity. Furthermore, the results reveal a decline for all methods in the percentage of images for which the estimated orientation is restricted to 1° as the orientation ambiguity increases. Nevertheless, our method and Shi *et al.* [28], which was recently accepted to ICCV2023, consistently maintain the majority of image orientations within a 3° margin from their ground truth values. Consequently, the translation estimation performance remains robust. In contrast, Shi and Li [26] encounter a notable drop in both translation and orientation estimation performance.

# 2 Performance with Increasing Amounts of Data as Supervision

Below, we analyze the performance of our method with  $\lambda = 0, 1$  and the stateof-the-art [26, 28], when different amounts of training data are employed. The results are illustrated in Fig. A.

For most models, except Shi and Li [26], we observe a consistent increase in performance on the same-area evaluation (Test-1) as the amount of training data increases. However, when it comes to the cross-area evaluation (Test-2),



**Fig. A:** Performance comparison between our method and the state-of-the-art on the KITTI dataset with different amounts of training data.

the two state-of-the-art methods, which require ground truth pose for supervision, exhibit a decline in performance when the training data exceeds 80%. This phenomenon suggests that our method avoids overfitting and holds the potential for further performance improvements with additional training data. Moreover, it's worth noting that our method doesn't necessitate GT labels for ground images during training, simplifying the process of large-scale data collection and reducing associated costs.

## 3 Semi-supervised Setting

Our method can be easily adapted to address the scenario when a small amount of training data with accurate pose labels is available by adding additional supervision to the network with this amount of data, e.g., using the training objective in Shi *et al.* [28].

In Tab. C, we show the performance of our method when fine-tuned with different portions of the current training dataset with accurate pose labels. Data amount = 0 indicates the original weakly supervised pipeline where no such data is available. The results show that fine-tuning leads to better performance when the data amount with accurate pose labels is over 50%. However, when the data amount is smaller, the fine-tuning leads to inferior performance. We suspect this is because the model overfits the limited training data with accurate labels, causing a loss of general inference ability on other data.

A properly designed training method should address this problem. However, apart from this semi-supervised scenario, we believe the original weaklysupervised application addressed in this paper is also of high importance. It 4 Shi et al.

**Table C:** The percentage of images whose lateral translation has been restricted to 1 meter (d = 1) to its ground truth value under the semi-supervised setting.

Data Amount (%)	0   0	0.2   0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
Lateral $(d=1)$	62.73 39	$9.49 \mid 44.61$	55.31	60.96	65.60	71.48	74.61	78.61	83.86

avoids additional effort to select the portion of accurate data, as most localization algorithms and noisy GPS sensors themselves do not provide a measure of whether its prediction is accurate or not. Thus, we leave this semi-supervised setting as a future work.

# 4 Sharing Feature Extractors or Not at Different Scenarios

We empirically found that the feature extractors for satellite and ground images captured by a pin-hole camera are shareable in the localization task. Tab. D presents the comparison of our method with shared or non-shared feature descriptors. In this comparison, the rotation estimator is kept the same and trained only on satellite images. From the results, it can be seen that sharing weights between ground and satellite images achieves better performance. A potential explanation might be that both the satellite images and ground images captured by a pin-hole camera map straight lines in the real world to straight lines in images and the viewpoint differences of the two view images are solved by a geometry projection module. This is similar to the task of multi-view stereo and image-based rendering, where the feature extractors for multi-view images are shared, and their differences are handled by Homography/geometry warping. Not sharing weights between the two branches increases the learning burden of the network, especially when supervision is not strong, resulting in inferior performance.

While for rotation and translation estimation, we found different feature extractors for different purposes achieve better performance. Tab. E illustrates the comparison results. This might be because good features for rotation and translation estimation are not identical. When re-using the feature extractors in rotation estimation for translation, we found the network converges slowly, and the performance on both rotation and translation estimation is poor, although better than the original coarse poses that we aim to refine.

#### 5 Performance on the Ford Dataset

We present the performance of our method on the Ford dataset in Tab. F. The results show our method achieves competitive performance with fully supervised SOTA, and sharing feature extractors between the ground and satellite branches performs better than not sharing, which is consistent with our observations on KITTI.

**Table D:** Sharing feature extractors or not between satellite and ground images captured by a pin-hole camera.

	Share?	$\begin{array}{c} \text{Lat} \\ d = 1 \uparrow \end{array}$	$\begin{array}{c} \text{Tes} \\ \text{eral} \\ d = 3 \uparrow a \end{array}$	st-1 (Sa Longit $l = 1 \uparrow$	$\begin{array}{l} \mathrm{ame-are} \\ \mathrm{udinal} \\ d=3\uparrow \end{array}$	a) Azir $\theta = 1 \uparrow$	$\theta = 3 \uparrow$	$\begin{array}{c} \text{Lat} \\ d = 1 \uparrow \end{array}$	$\begin{array}{c} \text{Te} \\ \text{eral} \\ d = 3 \uparrow \end{array}$	st-2 (C: Longit $d = 1 \uparrow$	ross-are udinal $d = 3 \uparrow$	a) Azir $\theta = 1 \uparrow$	$\begin{array}{l} \text{muth} \\ \theta = 3 \end{array} \uparrow$
Ours $(\lambda = 0)$	No Yes	$ \begin{array}{c}48.64\\59.58\end{array} $	$\begin{array}{c c} 77.37 \\ 85.74 \end{array}$	$9.97 \\ 11.37$	$\begin{array}{c} 25.63 \\ 31.94 \end{array}$	$99.66 \\ 99.66$	$\begin{array}{c} 100.00\\ 100.00 \end{array}$	$\begin{array}{c} 54.49 \\ 62.73 \end{array}$	$\begin{array}{c} 79.75 \\ 86.53 \end{array}$	$8.96 \\ 9.98$	$\begin{array}{c} 26.54 \\ 29.67 \end{array}$	$99.99 \\ 99.99$	$\begin{array}{c} 100.00\\ 100.00 \end{array}$
$0urs \\ (\lambda = 1)$	No Yes	$ \begin{array}{c} 62.81\\ 66.07 \end{array} $	93.00 94.22	$19.90 \\ 16.51$	55.53 49.96	$99.66 \\ 99.66$	$\begin{array}{c} 100.00\\ 100.00 \end{array}$	$\begin{array}{c} 62.81\\ 64.74 \end{array}$	84.77 86.18	$\begin{array}{c} 13.14\\ 11.81 \end{array}$	$36.89 \\ 34.77$	99.99 99.99	$\begin{array}{c} 100.00\\ 100.00 \end{array}$

Table E: Sharing feature extractors or not for rotation and translation estimation.

	Share?	$d = 1 \uparrow$	$ \begin{array}{c} \text{Te}\\ \text{eral}\\ d = 3 \uparrow \end{array} $	st-1 (Sa Longit $d = 1 \uparrow$	ame-are udinal $d = 3 \uparrow$	a) Azir $\theta = 1 \uparrow$	$\begin{array}{l} \text{nuth} \\ \theta = 3 \end{array} \uparrow$	$\begin{array}{c} \text{Lat} \\ d = 1 \uparrow \end{array}$	$ \begin{array}{c} \text{Te} \\ \text{eral} \\ d = 3 \uparrow \end{array} $	st-2 (C Longit $d = 1 \uparrow$	ross-are udinal $d = 3 \uparrow$	a) Azir $\theta = 1 \uparrow$	$\begin{array}{l} \text{nuth} \\ \theta = 3 \uparrow \end{array}$
$\begin{array}{c} \mathbf{Ours} \\ (\lambda = 0) \end{array}$	${ m Yes}$ <b>No</b>	33.77 59.58	$\begin{array}{c c} 74.66 \\ 85.74 \end{array}$	$9.17 \\ 11.37$	$\begin{array}{c} 26.13\\ 31.94 \end{array}$	$\begin{array}{c} 11.69 \\ 99.66 \end{array}$	$\begin{array}{c} 34.64 \\ 100.00 \end{array}$	$32.18 \\ 62.73$	$\begin{array}{c} 71.92 \\ 86.53 \end{array}$	$7.58 \\ 9.98$	$\begin{array}{c} 24.00\\ 29.67 \end{array}$	$\begin{array}{c} 12.64 \\ 99.99 \end{array}$	$\begin{array}{c} 37.24 \\ 100.00 \end{array}$
$\begin{array}{c c} \mathbf{Ours} \\ (\lambda = 1) \end{array}$	$\begin{array}{c} \mathrm{Yes} \\ \mathbf{No} \end{array}$	$ \begin{array}{c} 35.62\\ 66.07 \end{array}$	81.69 94.22	$\begin{array}{c} 10.68\\ 16.51 \end{array}$	$\begin{array}{c} 31.78\\ 49.96 \end{array}$	$\begin{array}{c} 10.20\\ 99.66\end{array}$	$31.35 \\ 100.00$	$\begin{vmatrix} 32.62 \\ 64.74 \end{vmatrix}$	73.72 86.18	8.87 11.81	$26.07 \\ 34.77$	$\begin{array}{c} 10.21\\ 99.99 \end{array}$	$\begin{array}{c} 31.62 \\ 100.00 \end{array}$

Table F: Average results comparison on Log1 and Log2 of the Ford Multi-AV dataset.

Algorithma		Lat	eral	Longit	udinal	Azir	nuth
Algorith	ins	$d = 1 \uparrow$	$d = 3 \uparrow$	$d = 1 \uparrow$	$d = 3 \uparrow$	$\theta = 1 \uparrow$	$\theta = 3 \uparrow$
Shi and Li	[26]*	37.29	69.33	5.22	15.94	16.87	44.14
Shi et al. $[28]^*$		63.77	82.53	19.45	34.17	61.49	92.62
Song et al.	[33]*	46.13	-	11.97	-	58.64	-
$\mathbf{O}_{\mathbf{u}\mathbf{r}\mathbf{c}}\left(1-0\right)$	share	65.47	93.50	8.00	20.11	46.63	96.86
Ours $(\lambda = 0)$	not share	14.21	69.62	5.13	13.40	46.63	96.86
Ours $(1 - 1)$	share	78.87	97.96	16.42	37.26	46.63	96.86
Ours $(\lambda = 1)$	not share	58.19	96.83	10.14	25.33	46.63	96.86

# 6 Visual Explanation of the Spatial Correlation Process

The spatial correlation process is illiustrated in Fig. B. We first center crop the synthesized overhead view feature map depicted in Fig. 5 (b) of the main paper and make its coverage around  $40 \times 40$  m<sup>2</sup>, as we consider scene contents within 20m to the camera location is the most important for the localization purpose. Then, we adopt the cropped overhead view feature map as the correlation kernel, similar to the yellow kernel in Fig. B, and the reference satellite image as the correlation input, indicated by the green grid in Fig. B, and apply inner product between the input and the kernal. The output, indicated by the pink grid in Fig. B, is the location probability map of the ground camera with respect to the satellite image.

In practice, the coverage of the reference satellite map and the kernel is engineered to make the coverage of the convolution output slightly larger than the location search space of the ground camera. In this example, the coverage of the satellite map is around  $100 \times 100$  m<sup>2</sup>, and that of the convolution output (location probability map) is about  $60 \times 60$  m<sup>2</sup>.

6 Shi et al.



Fig. B: The spatial correlation process. We compute the inner product between the reference satellite feature map (input) and synthesized overhead view feature map (kernel) from the ground image across all possible locations. This figure is from https://giphy.com/gifs/blog-daniel-keypoints-i4NjAwytgIRDW

Table G: Comparison between projecting features (Feat.) and images (Imgs).

d =	$\begin{array}{c c} \text{Test-1} \\ \text{Lateral} & \text{Lon} \\ 1 \uparrow d = 3 \uparrow d = \end{array}$	$egin{array}{llllllllllllllllllllllllllllllllllll$	$\begin{array}{c c} \text{muth} & \text{Lat} \\ \uparrow \theta = 3 \uparrow d = 1 \uparrow \end{array}$	Test-2 (C eral Longit $d = 3 \uparrow d = 1 \uparrow$	ross-area) audinal   Azimuth $d = 3 \uparrow \theta = 1 \uparrow \theta = 3 \uparrow$
$\begin{array}{c c} \mathbf{Ours} &  \operatorname{Imgs}  & 45\\ (\lambda = 0) &  \mathbf{Feat}  & 59 \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c c c} 7 & 21.39 & 99.92 \\ 37 & 31.94 & 99.66 \end{array}$	$\begin{array}{c c c} 100.00 & 48.24 \\ 100.00 & 62.73 \end{array}$	$\begin{array}{c c} 79.16 & 6.95 \\ 86.53 & 9.98 \end{array}$	20.60 100.00 100.00 29.67 99.99 100.00
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c c c} 100.00 & 58.53 \\ 100.00 & 64.74 \end{array}$	87.34 11.24 86.18 11.81	33.23100.00100.0034.7799.99100.00

## 7 Stage-wise Training VS. End-to-end Training

We present the comparison between end-to-end and stage-wise training in Tab. H. We observe that, at the beginning of end-to-end training, the rotation estimation performance for the query image is poor (near random), which negatively impacts the translation estimation performance. At the same time, allowing the signal of translation estimation loss to propagate back to the rotation estimator negatively affects its performance. Furthermore, end-to-end training requires large GPU memory and, thus, a smaller batch size, which is detrimental to metric learning performance.

**Table H:** Comparison with end-to-end training on KITTI dataset (Ours with  $\lambda = 0$ )

	•											
Training		Test-1 (Same-area	a)		Test-2 (Cross-area	a)						
Annnaach	Lateral	Longitudinal	Azimuth	Lateral	Longitudinal	Azimuth						
Approach	d = 1	$\bar{d} = 1$	$\theta = 1^{\circ}$	d = 1	$\overline{d} = 1$	$\theta = 1^{\circ}$						
End-to-end	32.83	8.22	10.15	32.53	7.76	10.10						
Stage-wise	59.58	11.37	99.66	62.73	9.98	99.99						

# 8 Angle ambiguity at 0°& 360°, why not rotation matrix?

Since we have prior knowledge of the coarse orientation, the angle ambiguity can be restricted to be smaller than  $360^{\circ}$ . Thus, the angle ambiguity at  $0^{\circ}$  and  $360^{\circ}$  can be avoided. A rotation matrix over-parameterizes the 1-DoF rotation

angle in the cross-view image matching, thus leading to inferior performance, as shown in the 1st row of Tab. I. We further make the pose regressor output sine and cosine of the angle (2-DoF output). The results in the 2nd row of Tab. I are promising but still inferior to our original parameterization (last row).

**Table I:** Additional ablation study on rotation estimation on the KITTI dataset.

Detetion Dependent or institut	Test-1 (S	Same-area)	Test-2 (Cross-area)			
Rotation Farameterization	$\theta = 1^{\circ}$	$\theta = 3^{\circ}$	$\theta = 1^{\circ}$	$\theta = 3^{\circ}$		
Rotation Matrix	5.05	15.44	5.07	15.51		
Sin(angle) & Cos(angle)	90.01	100.00	86.52	100.00		
Angle	99.66	100.00	99.99	100.00		

#### 9 Comparison Between Projecting Features and Images

In this paper, we follow the general practice of projecting features instead of images [32]. This is because when projecting ground images to an overhead view by assuming ground plane Homography, the pixels for scene objects above the ground plane are incorrectly projected to the overhead view and thus suffer distortion. In this way, the scene information of these objects will be lost in the projected image, resulting in inferior localization performance.

In contrast, features have a larger field of view of the original image and encode higher-level semantic information about the scene. For example, the building roots also have a semantic meaning of "building". It can be mapped to the building roof in the overhead view, which shares the same semantic information as the building root. Thus, projecting features instead of the original images can tolerate the errors in the overhead-view projection by the ground plane Homography to some extent. We illustrate the experimental comparison between projecting features and images in Tab. G. Not surprisingly, projecting features achieves better performance.

#### 10 Model Size and Evaluation Speed Comparison

We present the model size and evaluation speed comparison with two recent state-of-the-art, whose models and evaluation scripts have been released, in Tab. J. All of them are evaluated on an RTX 3090 GPU. It can be seen that our method achieves the fastest evaluation speed with a relatively small model size.

#### 11 Limitations

Although our self-supervised learning approach has achieved promising results, it has a few limitations.

(i) First, as explained previously, our self-supervised training strategy for rotation estimation is only suitable for ground images captured by a pin-hole

Model Size					Evaluation Speed					
Shi and Li [26]   Shi et al. [28]   Ours						Shi and Li $\left[ 26\right]$		Shi at al. $\left[28\right]$		Ours
20.2 M		29.1 M		$20.6~{\rm M}$		$500 \mathrm{~ms}$		$200 \mathrm{\ ms}$		$47~\mathrm{ms}$

Table J: Model size and evaluation speed comparison on the KITTI dataset.

camera. Due to the significant domain differences between panoramic and satellite images, it cannot be applied to estimating a spherical camera's orientation.

(ii) Second, our deep metric learning supervision strategy computes the spatial correlation between each query image and several satellite images. To save GPU memory and enable a reasonable batch size for metric learning, we use the feature level of a quarter of the original image size for the translation estimation. This actually sacrifices localization accuracy to some extent.

(iii) For the same reason, we cannot adopt a more sophisticated overheadview feature synthesis method because it will consume significant memory, thus sacrificing batch size, and the weak supervision limits the learning ability of a powerful overhead-view synthesis module with complex designs.

(vi) Finally, similar to all the ground-to-satellite localization networks where a single camera is used for query, our method suffers poor localization performance along the longitudinal direction, as shown in Fig. C. This can potentially be addressed using a video or multi-camera setup for the query.

We leave these unsolved problems as our future work and encourage the community to pay attention to them.



Fig. C: Ambiguity along the longitudinal (driving direction) pose estimation.

#### 12 Potential Negative Impact

This paper introduces a novel approach for self-supervised ground-to-satellite image registration. The objective of this approach is to enhance the accuracy of coarse camera pose estimates, such as those obtained from noisy GPS sensors, through ground-to-satellite image matching. However, it also raises concerns about privacy, particularly regarding the potential for individuals or sensitive locations to be identified and tracked without their consent. Unauthorized access to satellite imagery could enable surveillance activities or intrusions into personal privacy, raising ethical and legal implications.

We emphasize that the proposed method should be utilized in a manner that aligns with legal and ethical considerations. Careful implementation and adherence to privacy regulations and policies are crucial to ensure the ethical and responsible use of this approach.