

---

# Self-supervising Fine-grained Region Similarities for Large-scale Image Localization

---



**Yixiao Ge<sup>1</sup>, Haibo Wang<sup>3</sup>, Feng Zhu<sup>2</sup>, Rui Zhao<sup>2</sup>, Hongsheng Li<sup>1</sup>**

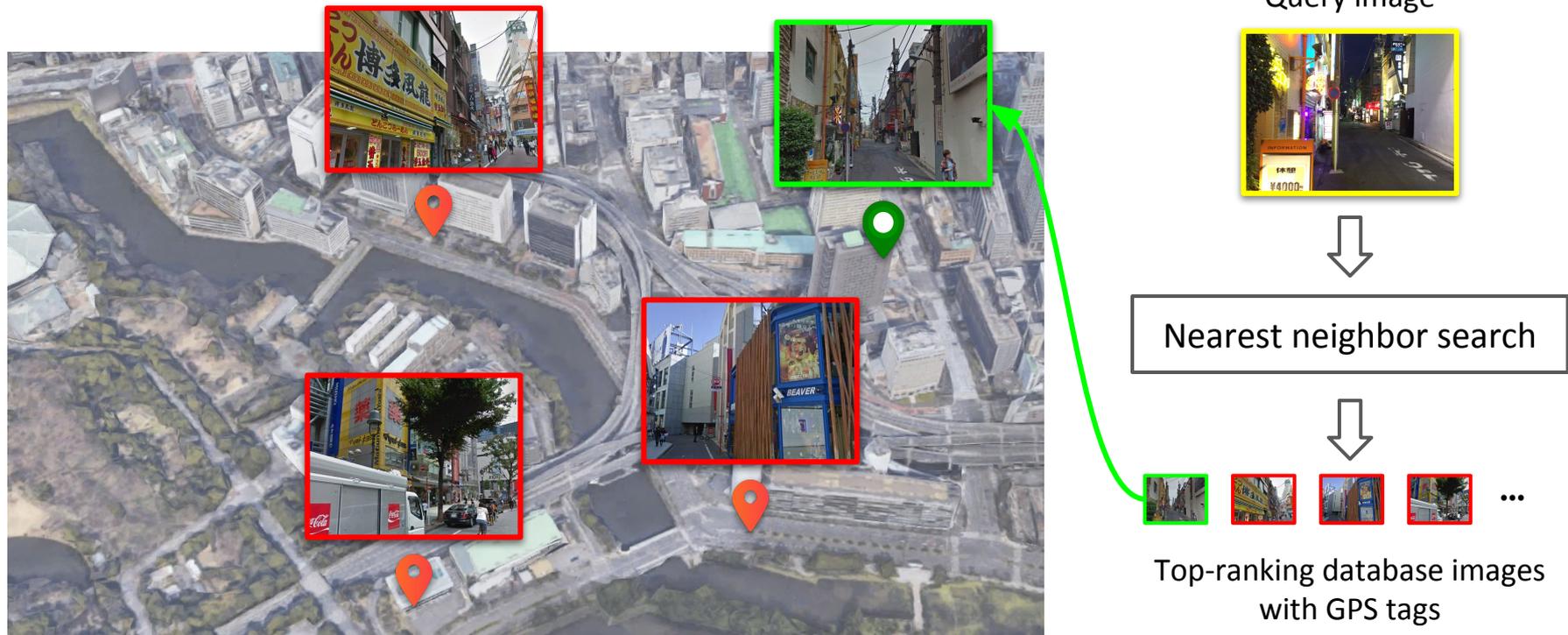
<sup>1</sup> The Chinese University of Hong Kong,

<sup>2</sup> SenseTime Research,

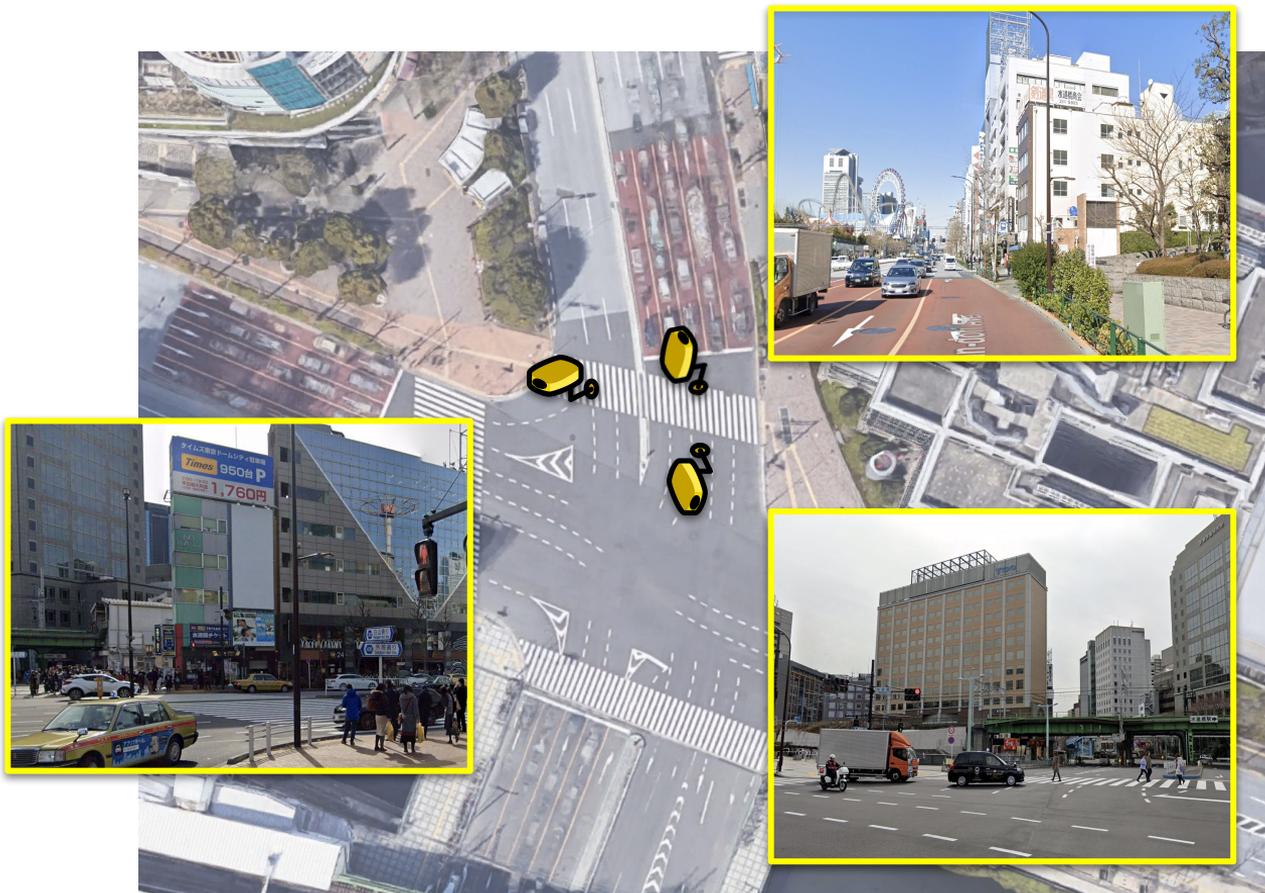
<sup>3</sup> China University of Mining and Technology



# Image Localization via Image Retrieval



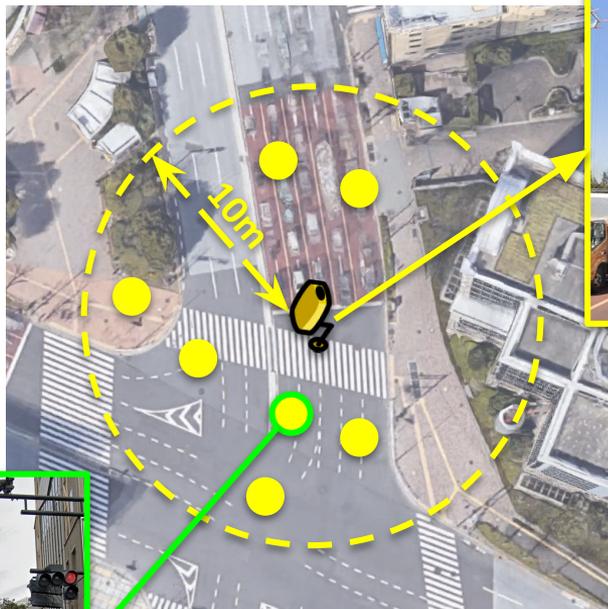
# Challenge #1: Noisy Positives by Weak GPS Labels



Geographically close-by images may not depict the same scene when facing different directions.

# Previous Solution: Train with Only the Easiest Positive

Potential positives filtered by GPS labels



Query image



Top-1 database image

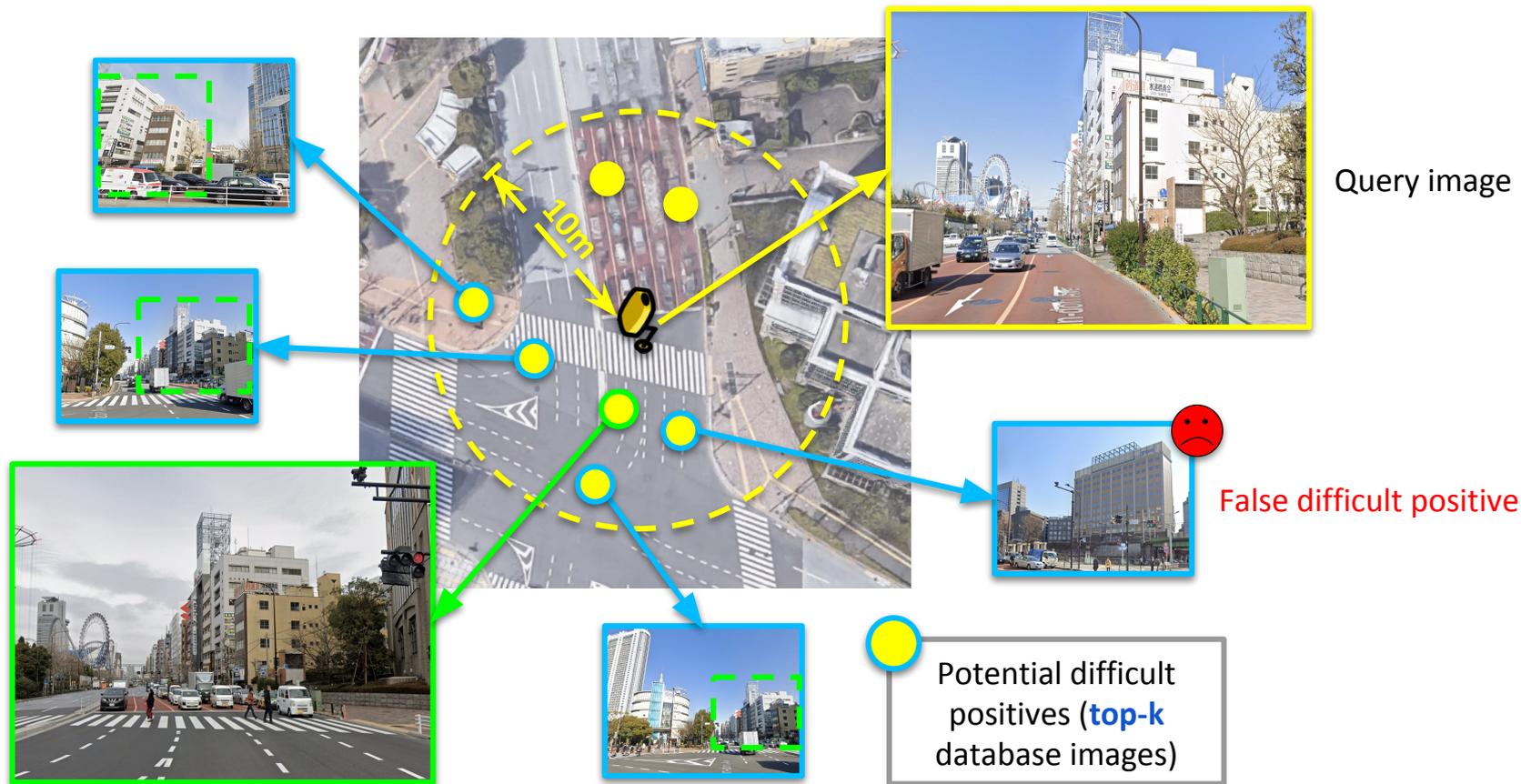


Forcing the queries to be closer to their already nearest neighbors results in a lack of robustness to varying conditions.



**Difficult positives are needed!**

# Motivation: Use Noisy Difficult Positives Properly



# Our Solution: Image Similarities as Soft Supervisions

Similarity label = 0.6



Similarity label = 0.5



Similarity label = 1.0



Similarity label = 0.3

Small similarity label for *true* difficult positive with *small overlapping regions*



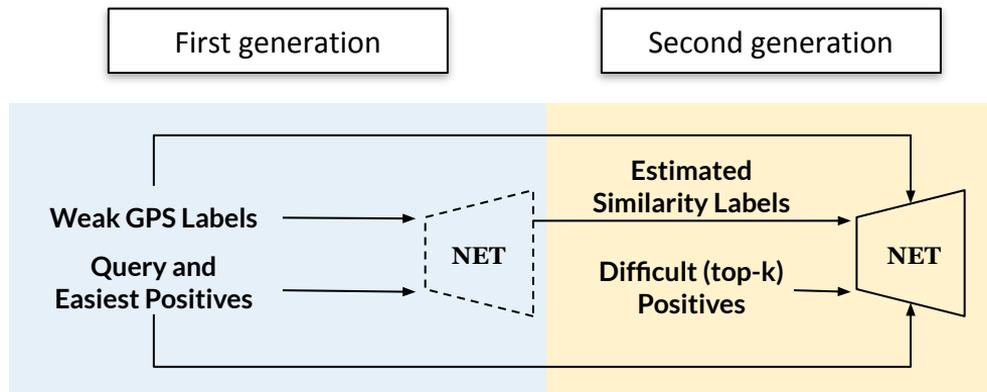
Query image



Similarity label = 0.1

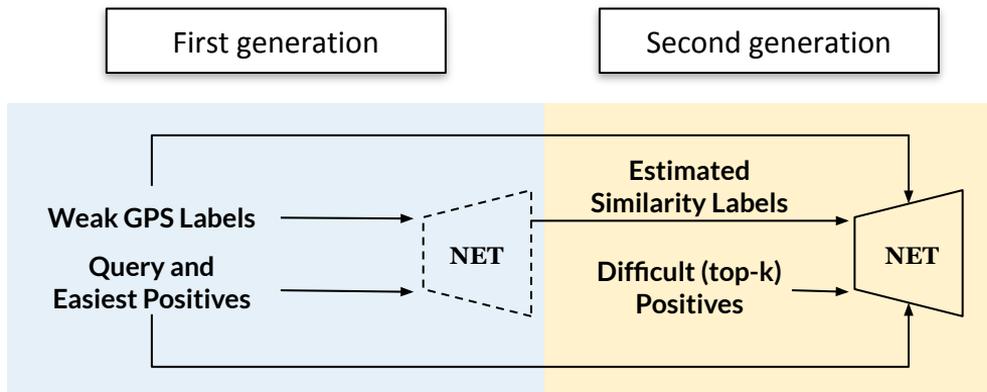
Small similarity label for *false* difficult positive

# Our Solution: Similarity Labels



The first generation's query-gallery similarities serve as the soft supervision for training the network in the second generation.

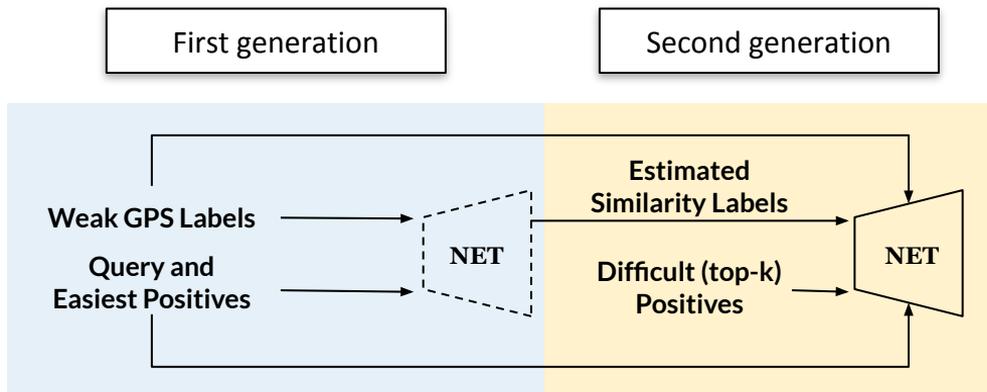
# Our Solution: Similarity Labels



**Similarity labels:**  $\mathcal{S}_{\theta_1}(q, p_1, \dots, p_k; \tau_1) = \text{softmax} \left( \left[ \frac{\langle f_{\theta_1}^q, f_{\theta_1}^{p_1} \rangle}{\tau_1}, \dots, \frac{\langle f_{\theta_1}^q, f_{\theta_1}^{p_k} \rangle}{\tau_1} \right]^\top \right)$

Query      Positive #1      Temperature for generation #1      Image similarity between query and positive #1      Parameters of the network in generation #1

# Our Solution: Similarity Labels

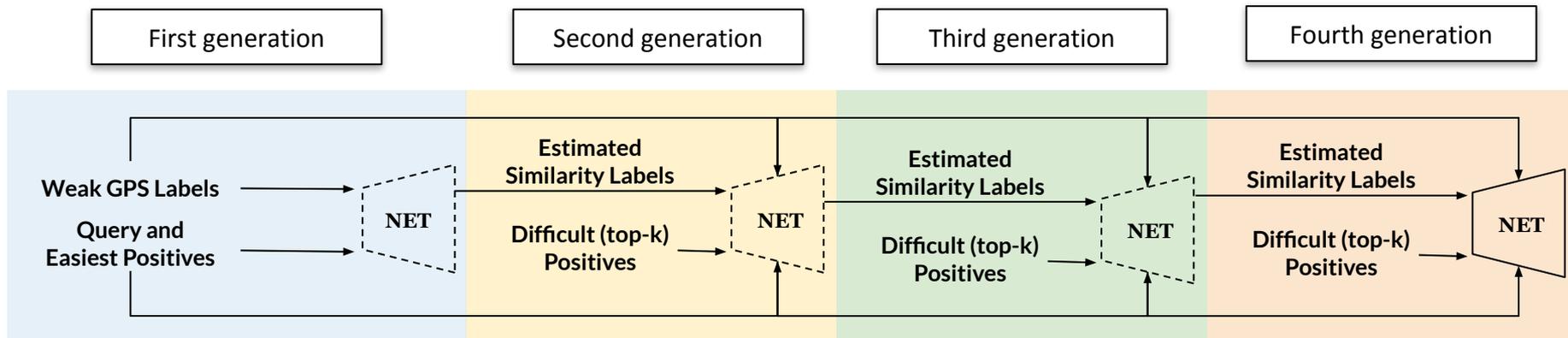


**Similarity labels:**  $\mathcal{S}_{\theta_1}(q, p_1, \dots, p_k; \tau_1) = \text{softmax} \left( \left[ \langle f_{\theta_1}^q, f_{\theta_1}^{p_1} \rangle / \tau_1, \dots, \langle f_{\theta_1}^q, f_{\theta_1}^{p_k} \rangle / \tau_1 \right]^\top \right)$

**Soft-label loss:**  $\mathcal{L}_{\text{soft}}(\theta_2) = \ell_{ce}(\mathcal{S}_{\theta_2}(q, p_1, \dots, p_k; 1), \mathcal{S}_{\theta_1}(q, p_1, \dots, p_k; \tau_1))$

↓ **cross-entropy loss**
↓ **similarity labels (learning targets)**  
estimated by the network in generation #1

# Our Solution: Self-enhanced Similarity Labels



The generated soft supervisions are gradually refined as the network generation progresses.

## Challenge #2: Lack of Region-level Supervisions

Only image-level labels



Query image



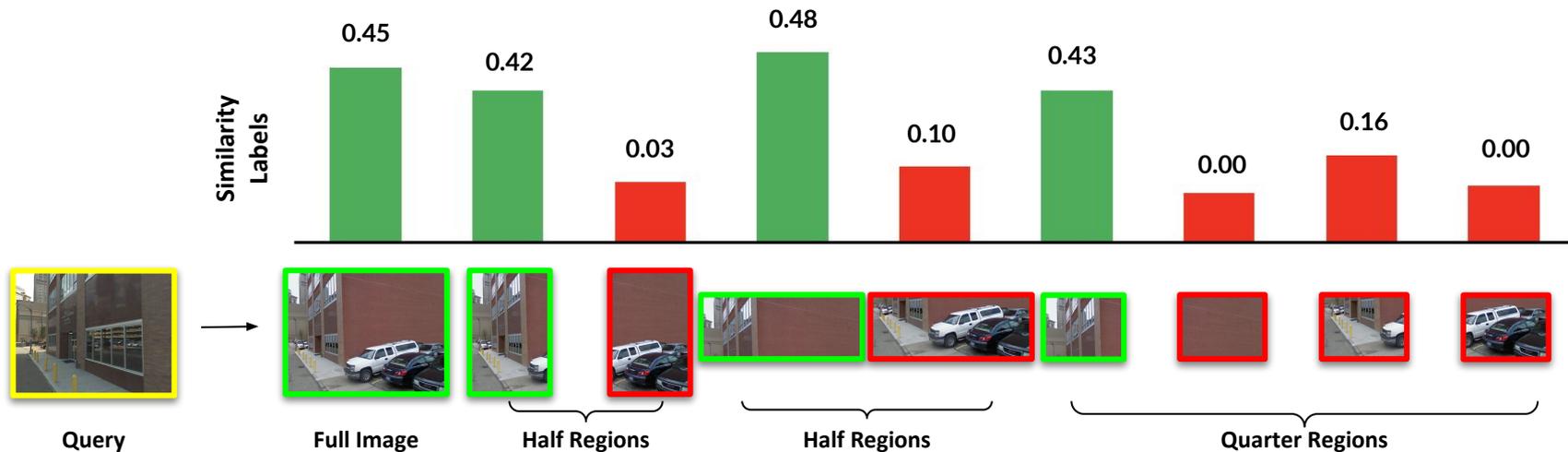
Positive sample

Ideal image-to-region labels



The correct image-level labels might not necessarily be the correct region-level labels.

# Our Solution: Image-to-region Similarities as Soft Supervisions



Provide fine-grained image-to-region similarities to enhance the learning of local features.

# Our Solution: Image-to-region Similarities as Soft Supervisions

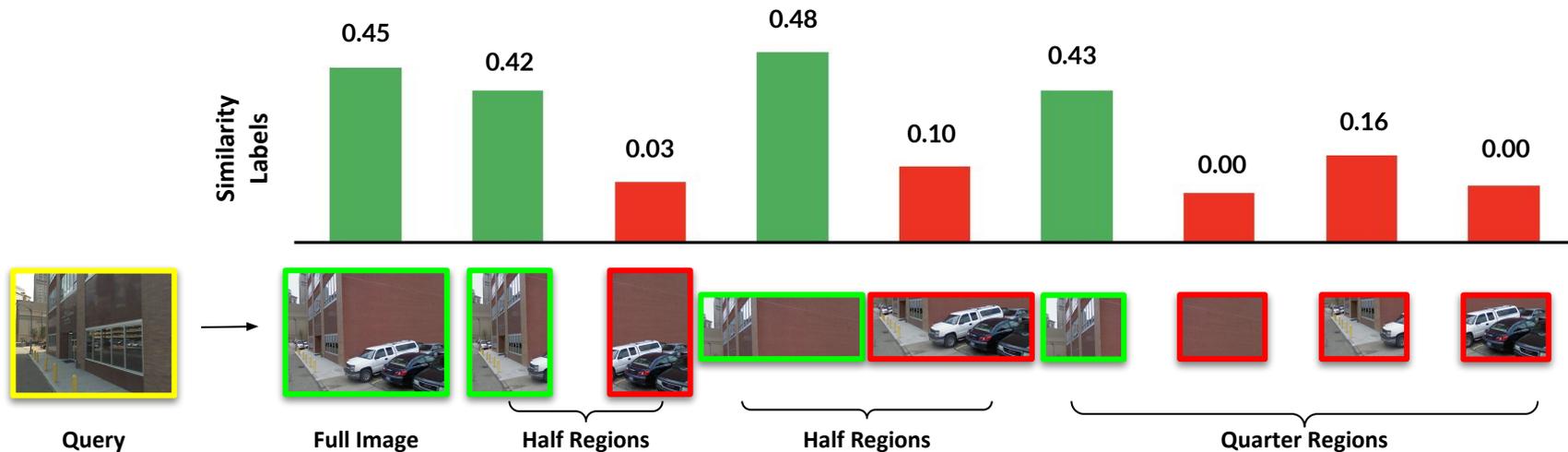


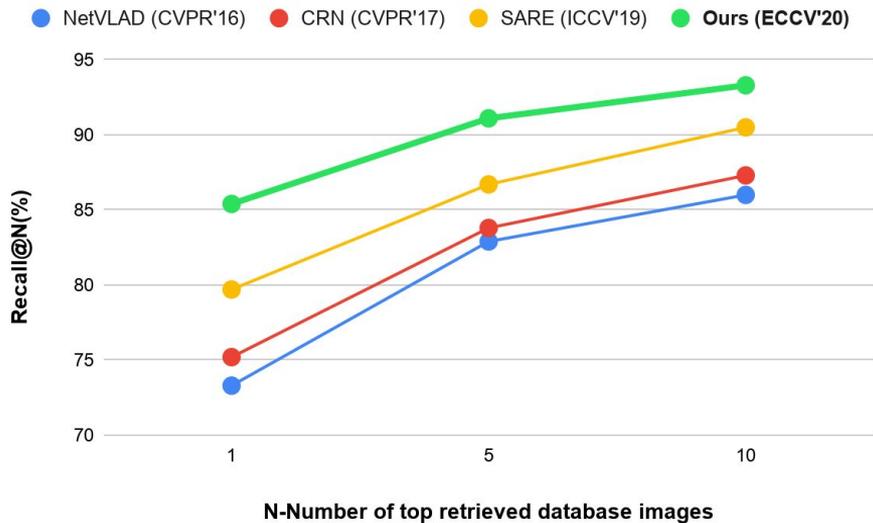
Image similarities between query and regions of positive #1

Fine-grained Similarity labels:

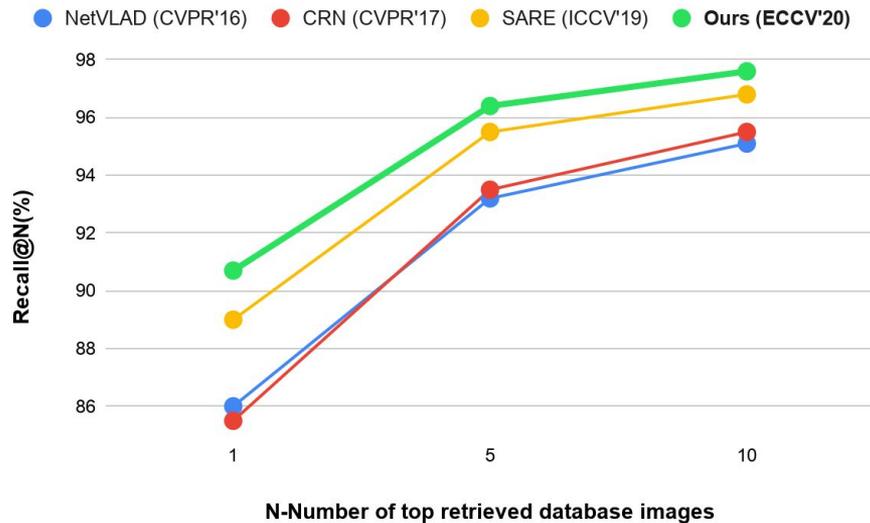
$$S_{\theta_{\omega}}^r(q, p_1, \dots, p_k; \tau_{\omega}) = \text{softmax} \left( \left[ \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{p_1} \rangle / \tau_{\omega}, \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{r_1^1} \rangle / \tau_{\omega}, \dots, \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{r_1^8} \rangle / \tau_{\omega}, \right. \right. \\ \left. \left. \dots, \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{p_k} \rangle / \tau_{\omega}, \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{r_k^1} \rangle / \tau_{\omega}, \dots, \langle f_{\theta_{\omega}}^q, f_{\theta_{\omega}}^{r_k^8} \rangle / \tau_{\omega} \right] \right)$$

# Performances on Image Localization Benchmarks

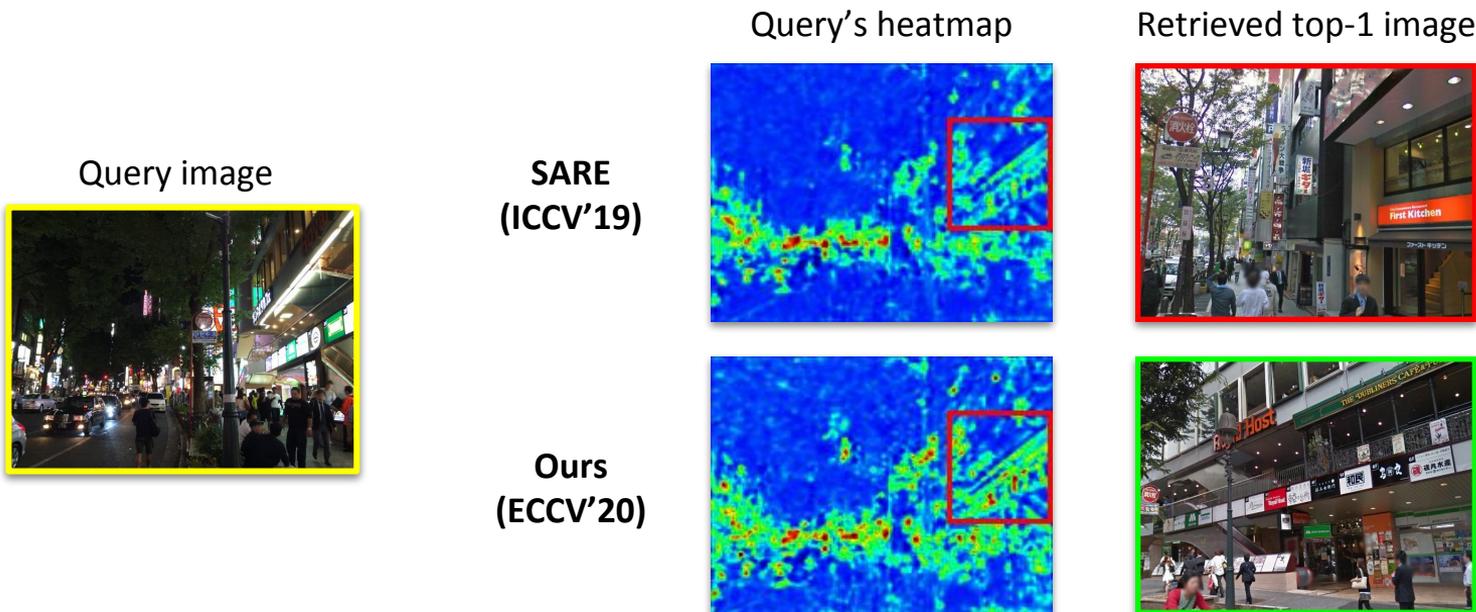
Tokyo 24/7



Pitts250k-test

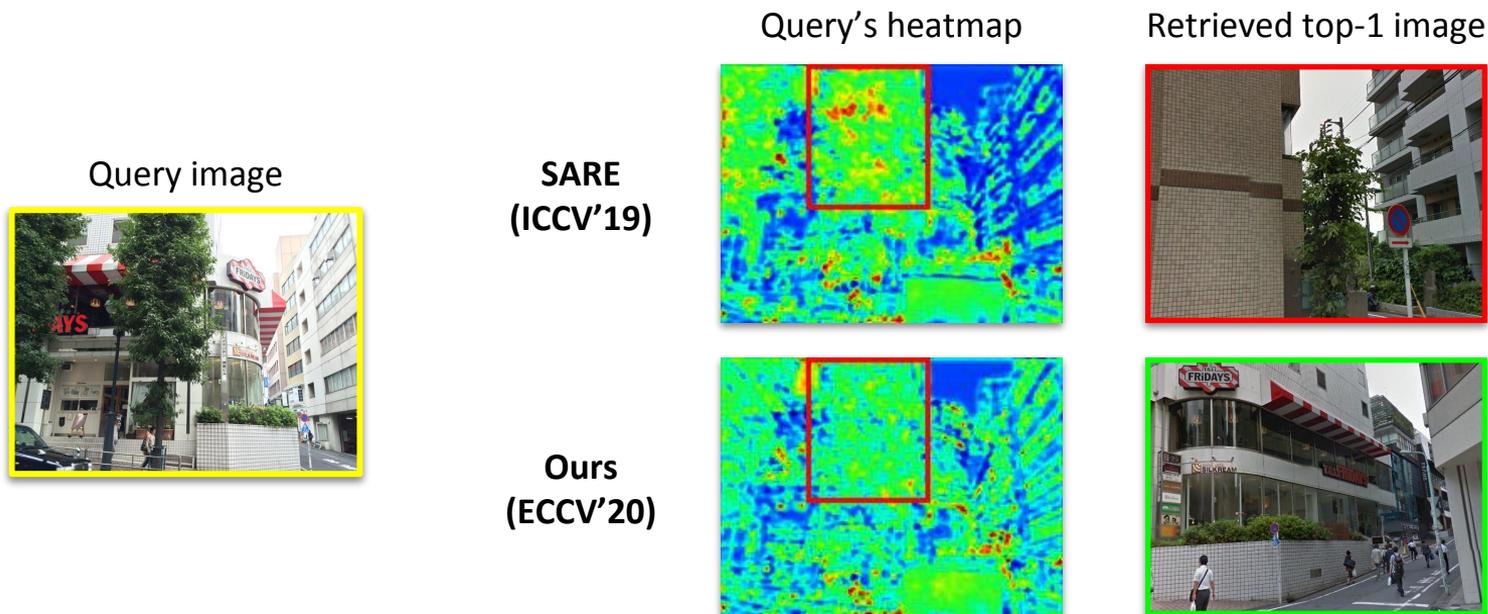


# Comparison with State-of-the-art (#1)



Our method pays more attention on the discriminative shop signs than SARE.

## Comparison with State-of-the-art (#2)



SARE incorrectly focuses on the trees, while our method learns to ignore such misleading regions.

# Self-supervising Fine-grained Region Similarities for Large-scale Image Localization

---

Yixiao Ge<sup>1</sup>, Haibo Wang<sup>3</sup>, Feng Zhu<sup>2</sup>, Rui Zhao<sup>2</sup>, Hongsheng Li<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong,

<sup>2</sup>SenseTime Research,

<sup>3</sup>China University of Mining and Technology

Code available at



<https://github.com/yxgeee/SFRS>