

# GeoRanker: Distance-Aware Ranking for Worldwide Image Geolocalization

**Pengyue Jia<sup>1,2</sup>, Seongheon Park<sup>2</sup>, Song Gao<sup>3</sup>, Xiangyu Zhao<sup>1</sup>, Sharon Li<sup>2</sup>,**

<sup>1</sup>Department of Data Science, City University of Hong Kong,

<sup>2</sup>Department of Computer Sciences, University of Wisconsin-Madison

<sup>3</sup>Department of Geography, University of Wisconsin-Madison

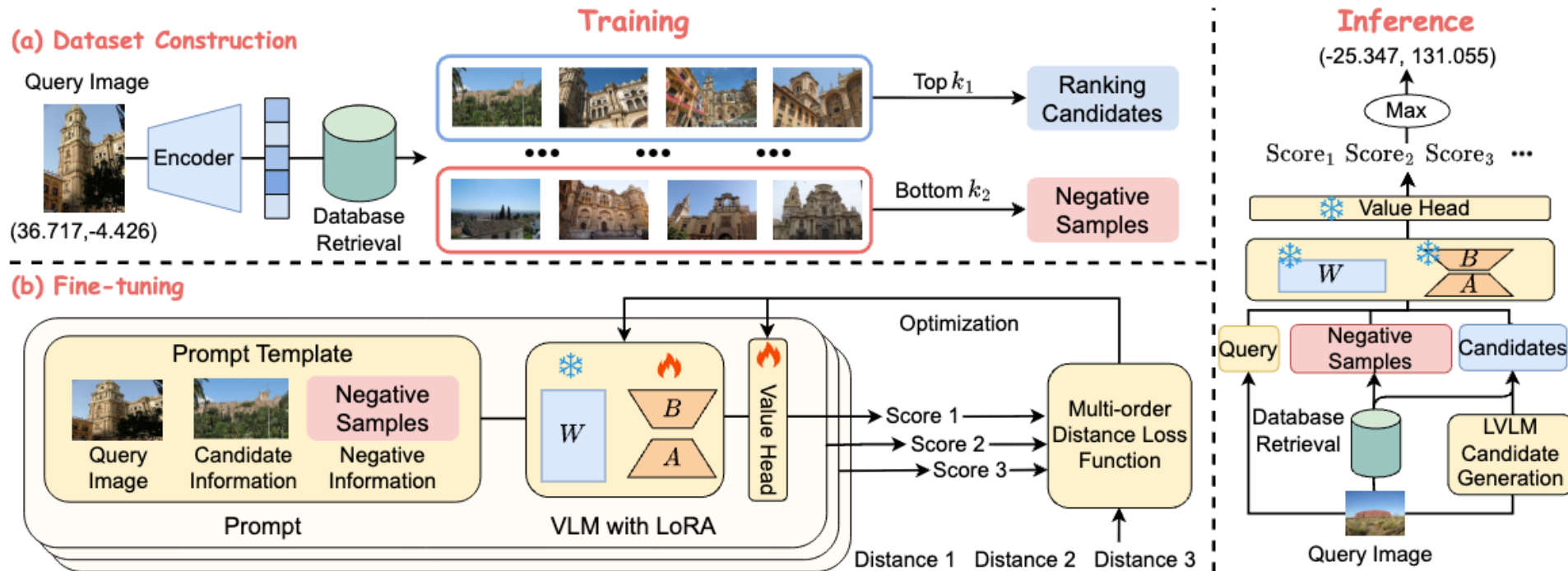
`jia.pengyue@my.cityu.edu.hk, sharonli@cs.wisc.edu`

## 01 Background & Motivation

## 02 Methodology

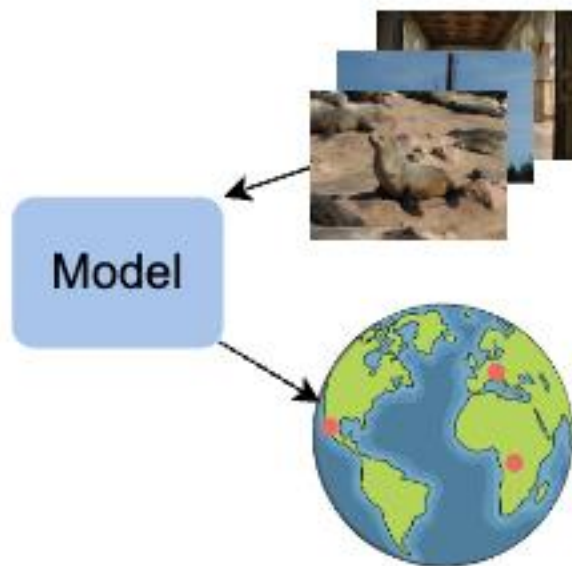
## 03 Experiments

## 04 Conclusion

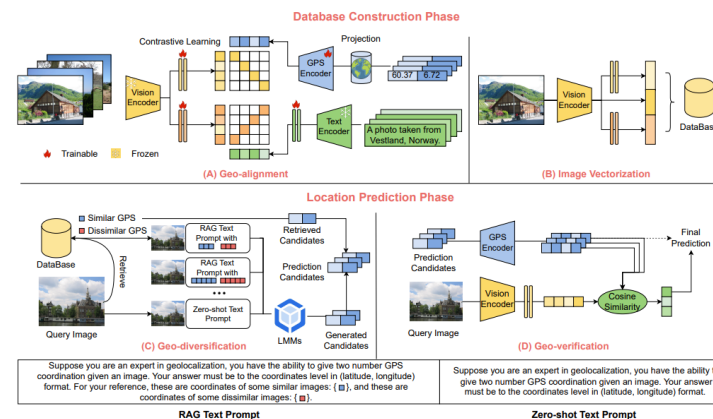


## Two-stage Pipeline is Widely Adopted

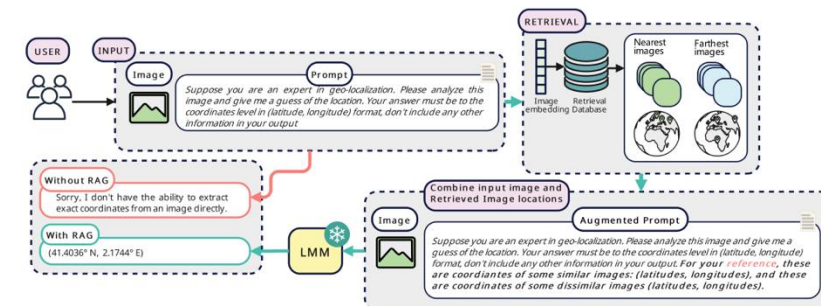
- Retrieve candidate locations from a global database
- Select the top match based on similarity scores



Worldwide Geolocalization



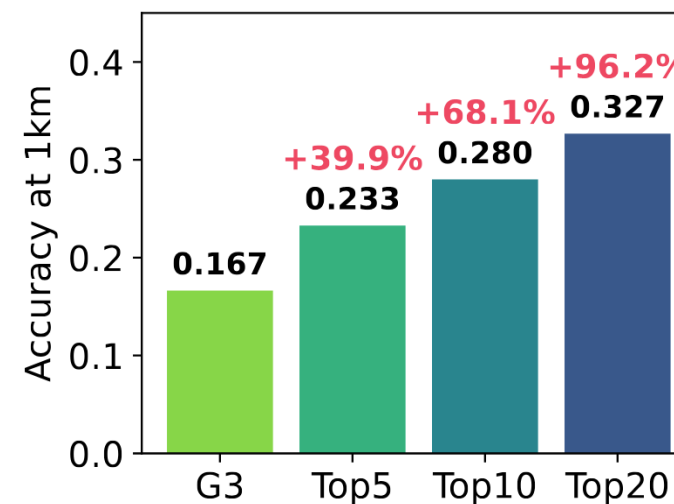
G3 [NeurIPS'24]



Img2Loc [SIGIR'24]

## Key Observation

Better candidates often exist within the top-k, but are not selected



## Root Cause

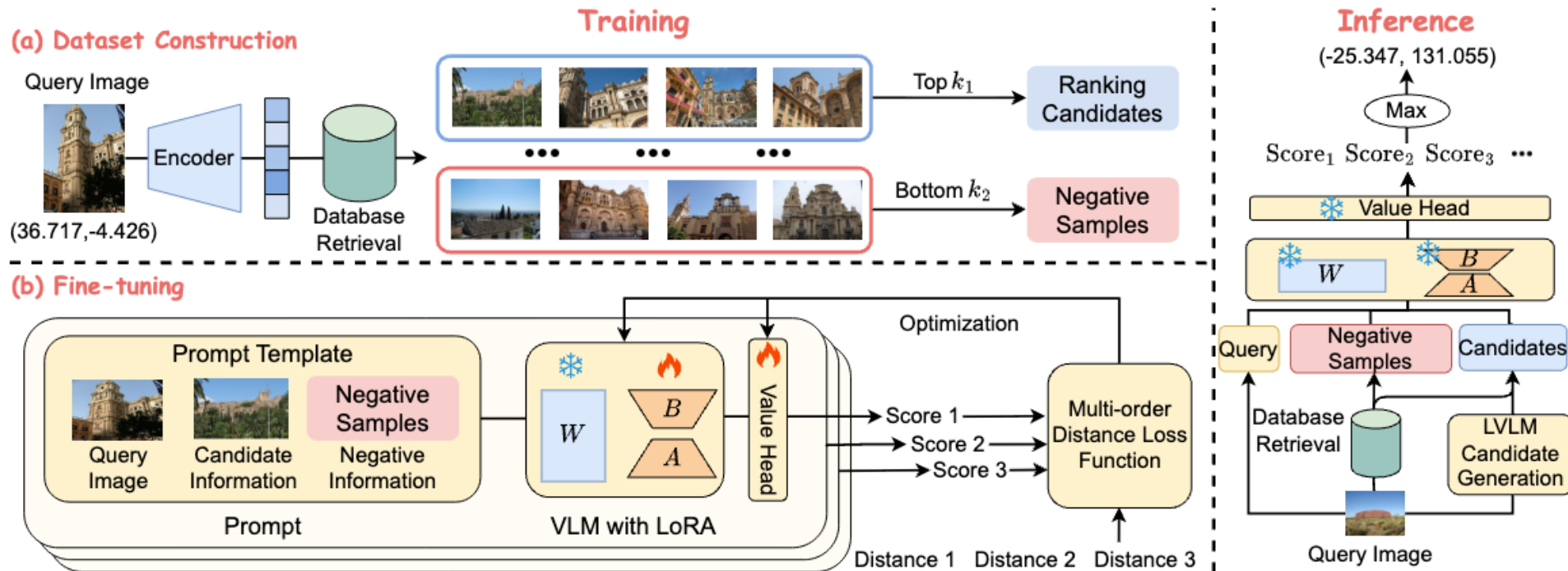
- Existing methods rely on simple similarity heuristics (e.g., cosine similarity of image embeddings)
- Existing training objectives primarily focus on point-wise similarity between individual images and locations, overlooking the rich spatial relationships among candidates.

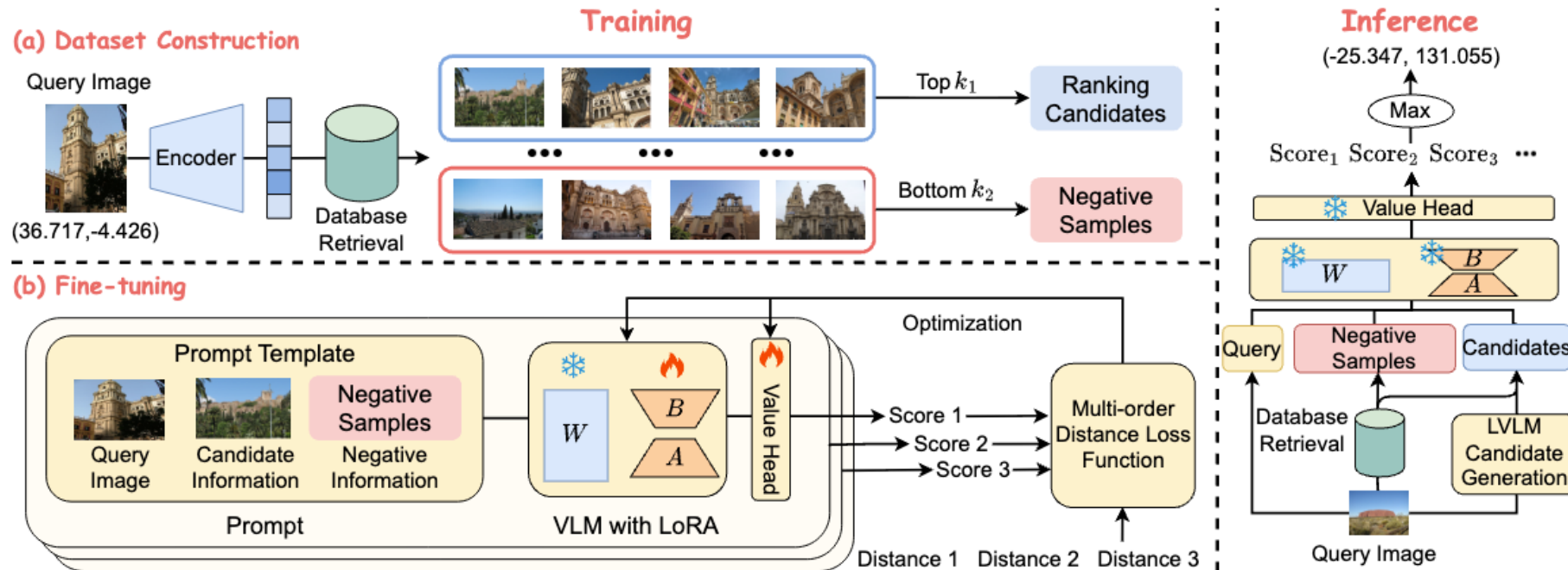
## 01 Background & Motivation

## 02 Methodology

## 03 Experiments

## 04 Conclusion

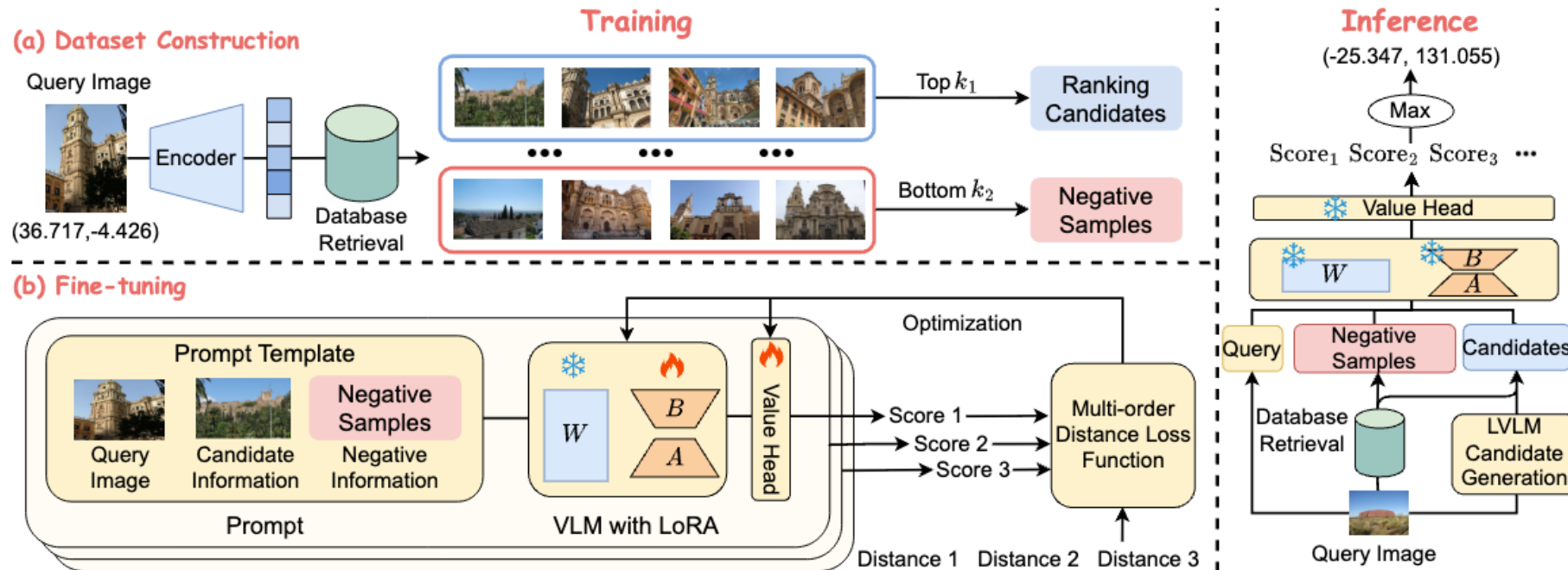




## GeoRanking Dataset Construction

- Candidate Encoding  $\mathbf{v}_{c_m} = \text{concat}(\text{Encoder}_c^{\text{gps}}(c_m^{\text{gps}}), \text{Encoder}_c^{\text{text}}(c_m^{\text{text}}), \text{Encoder}^{\text{img}}(c_m^{\text{img}}))$
- Random Sampling Query Image, Query Encoding
- Retrieving Top-N Candidates  $\mathbf{v}_q = \text{concat}(f_{\text{img} \rightarrow \text{gps}}(\text{Encoder}^{\text{img}}(q)), f_{\text{img} \rightarrow \text{text}}(\text{Encoder}^{\text{img}}(q)), \text{Encoder}^{\text{img}}(q))$



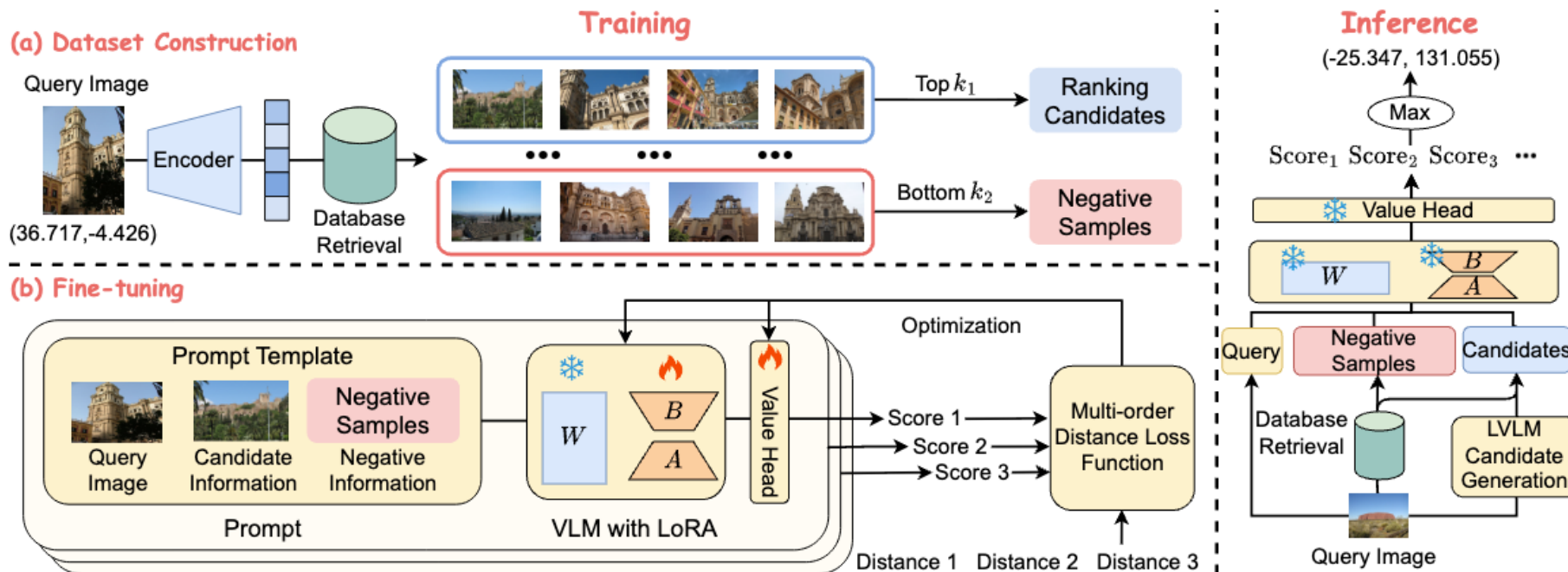


{query image} How far is this place from latitude: {candidate latitude}, longitude: {candidate longitude}, {candidate textual descriptions}, {candidate image}? Negative examples: {negative information}.

## GeoRanker

- Prompt
- Model Architecture: LoRA + Value Head

$$s = \mathbf{w}^\top \mathbf{h}_{\text{final}}, \quad \text{where } \mathbf{h}_{\text{final}} = \text{LVLM}(\mathbf{x})_{[-1]}$$



## GeoRanker: Optimization with Multi-Order Distance Objective

- First-order Distance Loss
- Second-order Distance Loss
- Joint Optimization

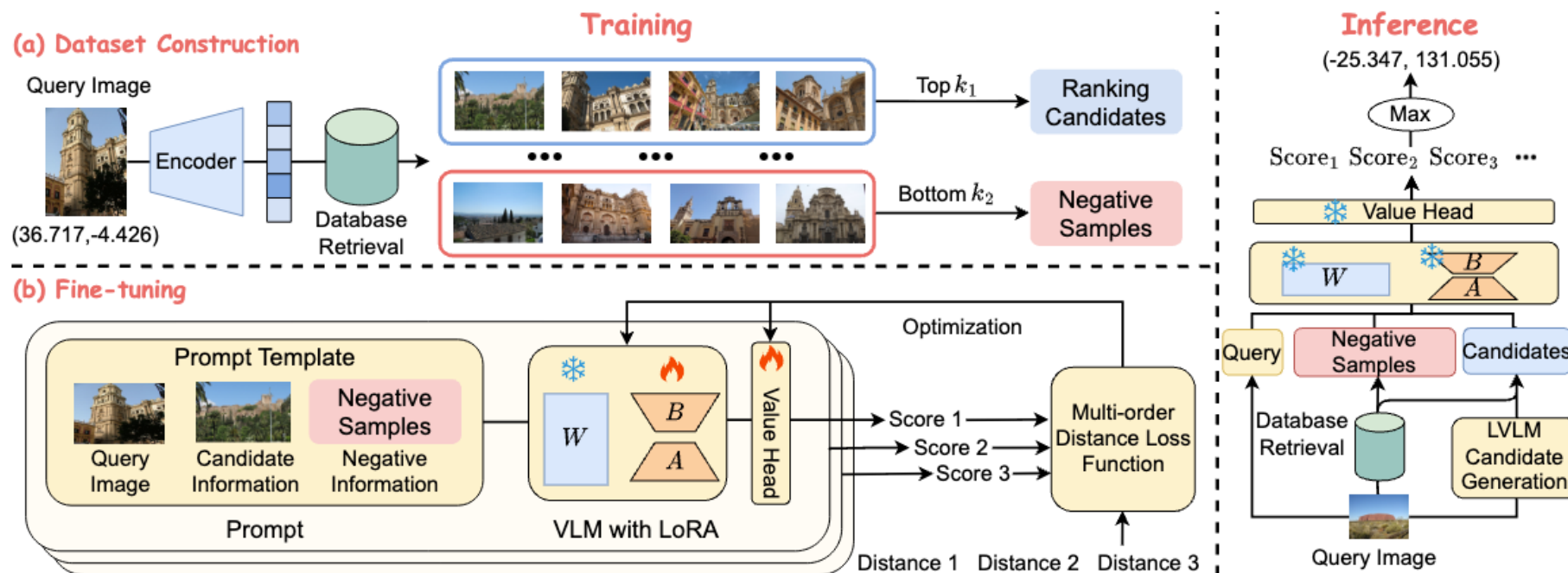
$$\Delta d_{i,j} = d_{\pi(i)} - d_{\pi(j)}, \quad \Delta s_{i,j} = s_{\pi(i)} - s_{\pi(j)}, \quad \text{for } 1 \leq i < j \leq k_1$$

$$\mathcal{L}_{\text{PL}}^{(2)} = -\frac{1}{K^{(2)}} \sum_{i=1}^{K^{(2)}} \log \frac{\exp(\Delta s_{(i)})}{\sum_{j=i}^P \exp(\Delta s_{(j)})}$$

$$\mathcal{L}_{\text{PL}}^{(1)} = -\frac{1}{K^{(1)}} \sum_{i=1}^{K^{(1)}} \log \frac{\exp(s_{\pi(i)})}{\sum_{j=i}^{k_1} \exp(s_{\pi(j)})}$$

$$\mathcal{L}_{\text{total}} = \lambda \cdot \mathcal{L}_{\text{PL}}^{(1)} + (1 - \lambda) \cdot \mathcal{L}_{\text{PL}}^{(2)}$$





## Inference

$$s_c = \text{GeoRanker}(q, c), \quad \forall c \in \mathcal{C}_r \cup \mathcal{C}_g$$

$$\hat{c} = \arg \max_{c \in \mathcal{C}_r \cup \mathcal{C}_g} s_c$$

## 01 Background & Motivation

## 02 Methodology

## 03 Experiments

## 04 Conclusion

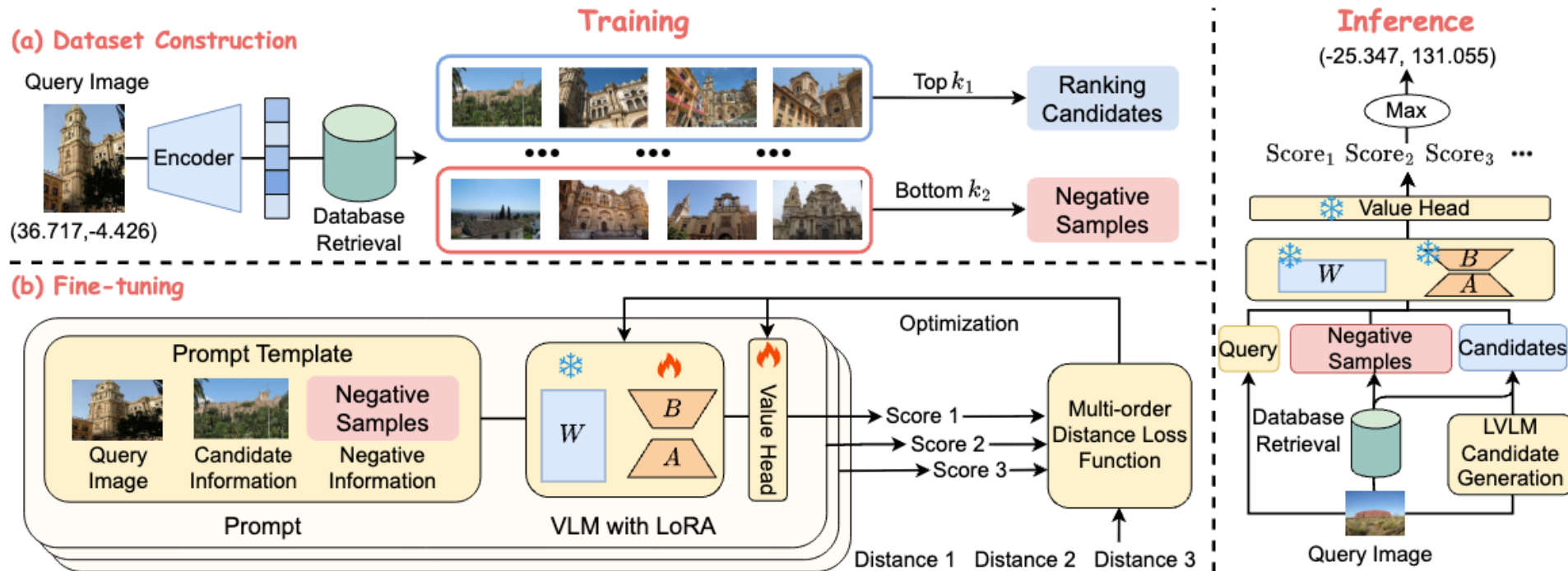


Table 1: **Main results** on IM2GPS3K and YFCC4K. For all metrics, higher is better. The best-performing results are highlighted in **bold**, while the second-best results are underlined.  $\Delta$  represents the relative improvement of our method over the best baseline.

Methods		IM2GPS3K					YFCC4K				
		Street 1km	City 25km	Region 200km	Country 750km	Continent 2500km	Street 1km	City 25km	Region 200km	Country 750km	Continent 2500km
[L]kNN, sigma=4 [1]	ICCV'17	7.2	19.4	26.9	38.9	55.9	2.3	5.7	11	23.5	42
PlaNet [24]	ECCV'16	8.5	24.8	34.3	48.4	64.6	5.6	14.3	22.2	36.4	55.8
CPlaNet [15]	ECCV'18	10.2	26.5	34.6	48.6	64.6	7.9	14.8	21.9	36.4	55.5
ISNs [55]	ECCV'18	10.5	28	36.6	49.7	66	6.5	16.2	23.8	37.4	55
Translocator [25]	ECCV'22	11.8	31.1	46.7	58.9	80.1	8.4	18.6	27	41.1	60.4
GeoDecoder [26]	ICCV'23	12.8	33.5	45.9	61	76.1	10.3	24.4	33.9	50	68.7
GeoCLIP [8]	NeurIPS'23	14.11	34.47	50.65	69.67	83.82	9.59	19.31	32.63	55	74.69
Img2Loc [10]	SIGIR'24	15.34	39.83	53.59	69.7	82.78	19.78	30.71	41.4	58.11	74.07
PIGEON [9]	CVPR'24	11.3	36.7	53.8	<u>72.4</u>	<u>85.3</u>	10.4	23.7	40.6	62.2	77.7
G3 [14]	NeurIPS'24	<u>16.65</u>	<u>40.94</u>	<u>55.56</u>	<u>71.24</u>	<u>84.68</u>	<u>23.99</u>	<u>35.89</u>	<u>46.98</u>	<u>64.26</u>	<u>78.15</u>
<b>GeoRanker</b>	<b>Ours</b>	<b>18.79</b>	<b>45.05</b>	<b>61.49</b>	<b>76.31</b>	<b>89.29</b>	<b>32.94</b>	<b>43.54</b>	<b>54.32</b>	<b>69.79</b>	<b>82.45</b>
Rel. Improvement	$\Delta$	$\uparrow 12.9\%$	$\uparrow 10.0\%$	$\uparrow 10.7\%$	$\uparrow 5.4\%$	$\uparrow 4.7\%$	$\uparrow 37.3\%$	$\uparrow 21.3\%$	$\uparrow 15.6\%$	$\uparrow 8.6\%$	$\uparrow 5.5\%$

- Superior Performance (37.3% improvement on YFCC4K in Street level)
- State-of-the-art across all datasets and metrics

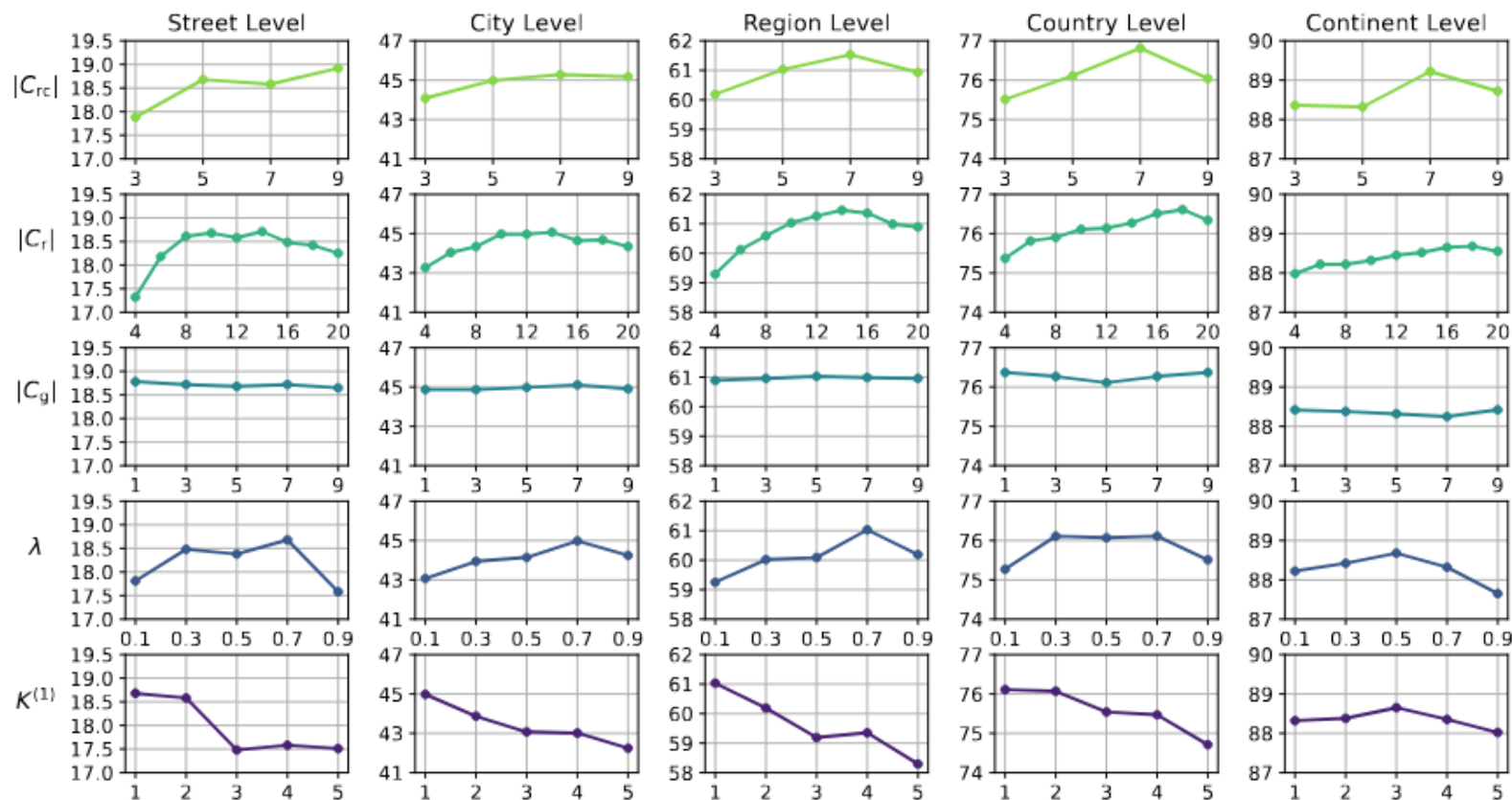
Methods	Street 1km	City 25km	Region 200km	Country 750km	Continent 2500km
w/o $\mathcal{L}_{PL}^{(2)}$	18.48	44.61	60.96	75.61	88.28
w/o $\mathcal{C}_{neg}$	17.35	44.51	60.82	76.37	88.28
w/o $\mathcal{C}_m^{text}$	18.02	43.91	60.19	<b>76.61</b>	88.62
w/o $\mathcal{C}_m^{img}$	15.58	41.77	59.15	75.40	88.35
w/o $\mathcal{C}_g$	18.21	43.47	59.69	75.47	88.75
Ours	<b>18.79</b>	<b>45.05</b>	<b>61.49</b>	76.31	<b>89.29</b>

## Ablation Study

Methods	IM2GPS3K				
	Street 1km	City 25km	Region 200km	Country 750km	Continent 2500km
Random	10.04	29.72	42.17	57.82	75.24
Top1	13.31	34.03	45.48	61.56	78.04
Prompting	16.62	40.21	54.55	70.07	83.24
Ours	<b>18.79</b>	<b>45.05</b>	<b>61.49</b>	<b>76.31</b>	<b>89.29</b>

## Comparison with Other Ranking Baselines

- All components contribute positively
- Removing any of the modality-aware prompt components leads to performance drops
- Without generated candidates underperforms GeoRanker
- GeoRanker is superior to Random, Top-1 Selection, and Prompting baselines.



- Impact of candidate scales in training and inference:  $C_{rc}$ ,  $C_r$ ,  $C_g$
- Impact of hyperparameters in multi-order distance objective:  $\lambda$ ,  $K^{(1)}$



Query Image



Top-5 Candidates without GeoReranker



870 KM > 69 KM > 0.44 KM < 596 KM > 440 KM

Top-5 Candidates with GeoReranker



0.44 KM < 69 KM < 440 KM < 596 KM < 870 KM

## Case Study

Parameter	Setting
GPU	NVIDIA L40S * 4
Training Time	16 hours / epoch
Total params	8,298,256,896
Trainable params	6,881,280 (0.083%)
Dataset Samples	100K
Batch Size	4
Batch Size per Device	1
Training GPU Memory Consumption	30 GB / GPU
VLM Backbone	Huggingface Qwen2-VL-7b-Instruct
Deepspeed	Stage 2

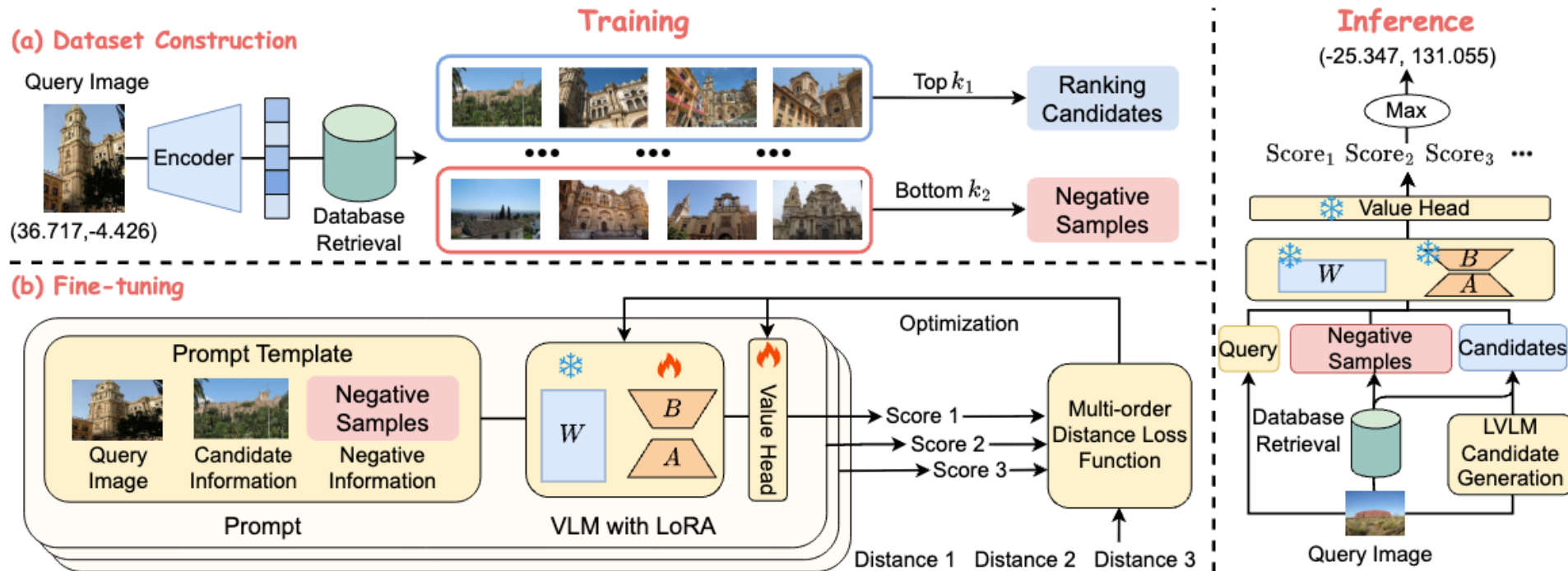
More Information on Training and Inference

## 01 Background & Motivation

## 02 Methodology

## 03 Experiments

## 04 Conclusion



- In this paper, we propose GeoRanker, a distance-aware ranking framework built upon LVLM.
- To enhance training, we introduce a novel multi-order distance loss that captures both absolute distances and relative spatial relationships among candidate locations.
- To support this framework, we construct GeoRanking, the first dataset specifically designed for spatial ranking tasks.
- Extensive experiments on IM2GPS3K and YFCC4K demonstrate the effectiveness of GeoRanker over baselines.



Code



MP16-Pro Dataset



CityU AML Lab



Pengyue's HomePage

# Thanks

JIA Pengyue

Applied Machine Learning Lab  
City University of Hong Kong  
[jia.pengyue@my.cityu.edu.hk](mailto:jia.pengyue@my.cityu.edu.hk)

Department of Computer Sciences  
University of Wisconsin - Madison  
[pjia7@wisc.edu](mailto:pjia7@wisc.edu)