



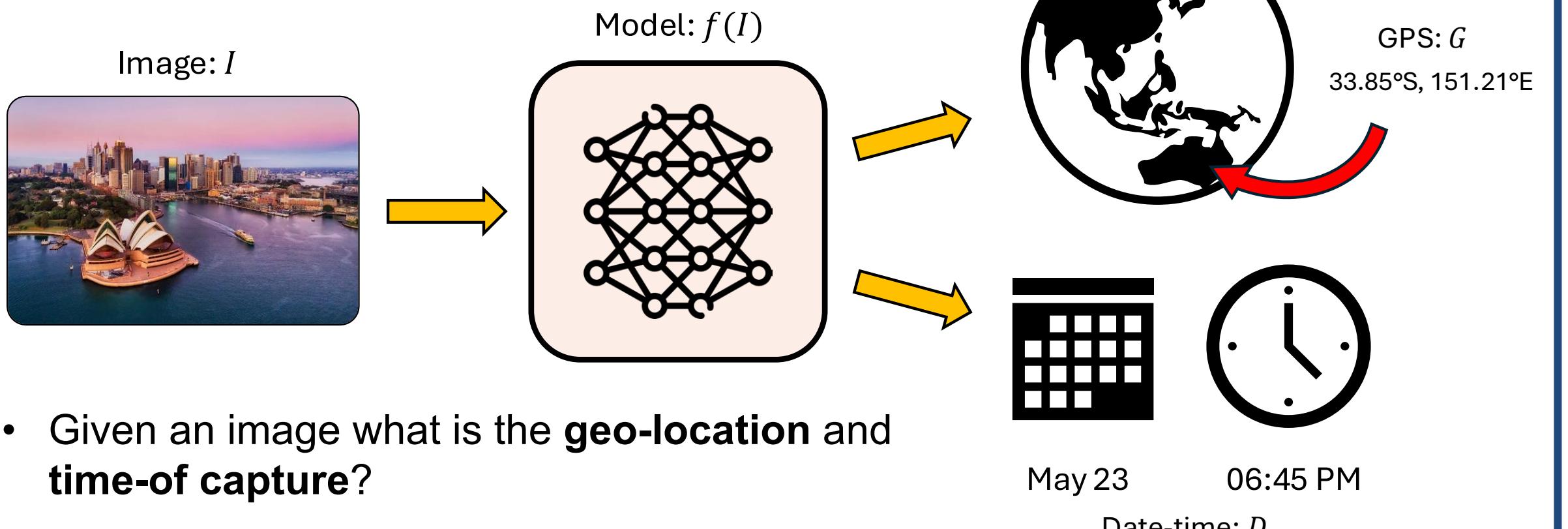
# GT-Loc: Unifying When and Where in Images Through a Joint Embedding Space

David Shatwell<sup>1</sup>, Ishan Dave<sup>2</sup>, Sirnam Swetha<sup>1</sup>, Mubarak Shah<sup>1</sup>

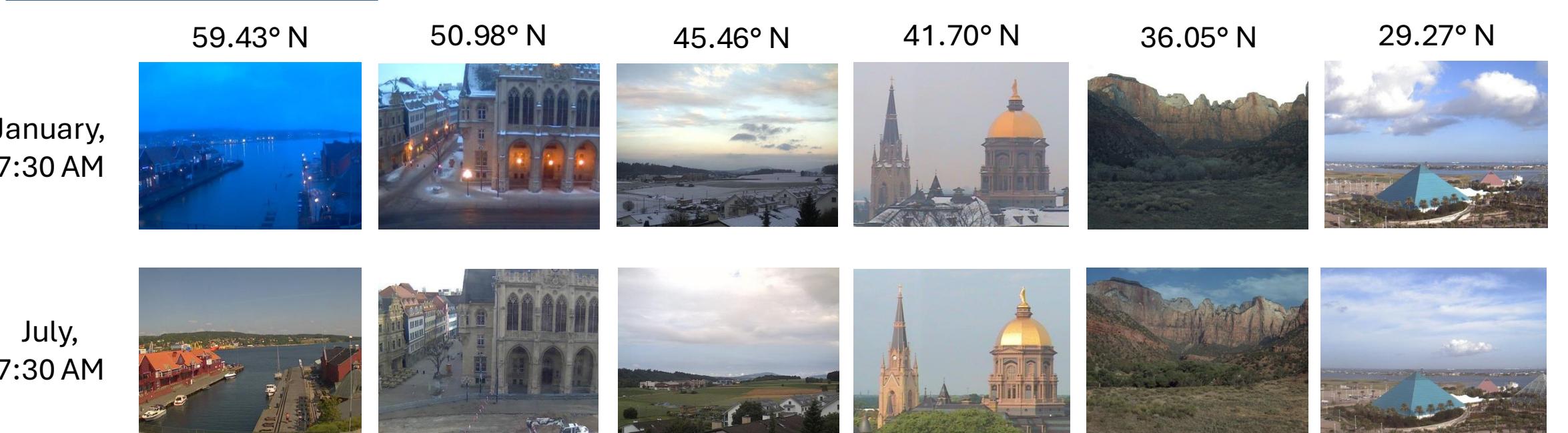
<sup>1</sup>University of Central Florida, <sup>2</sup>Adobe

★ ORAL

## GEO-TEMPORAL LOCALIZATION

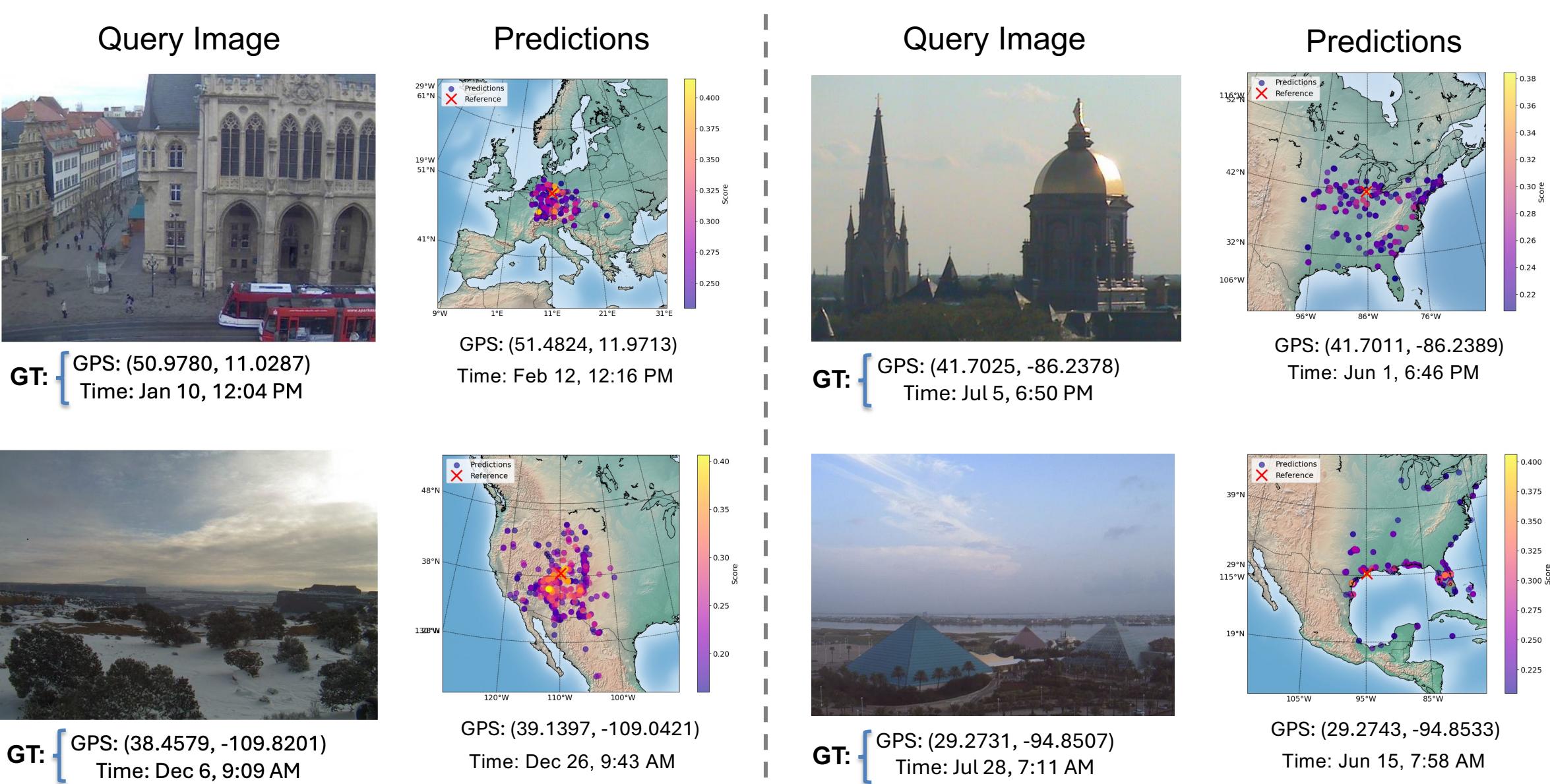


## CHALLENGES

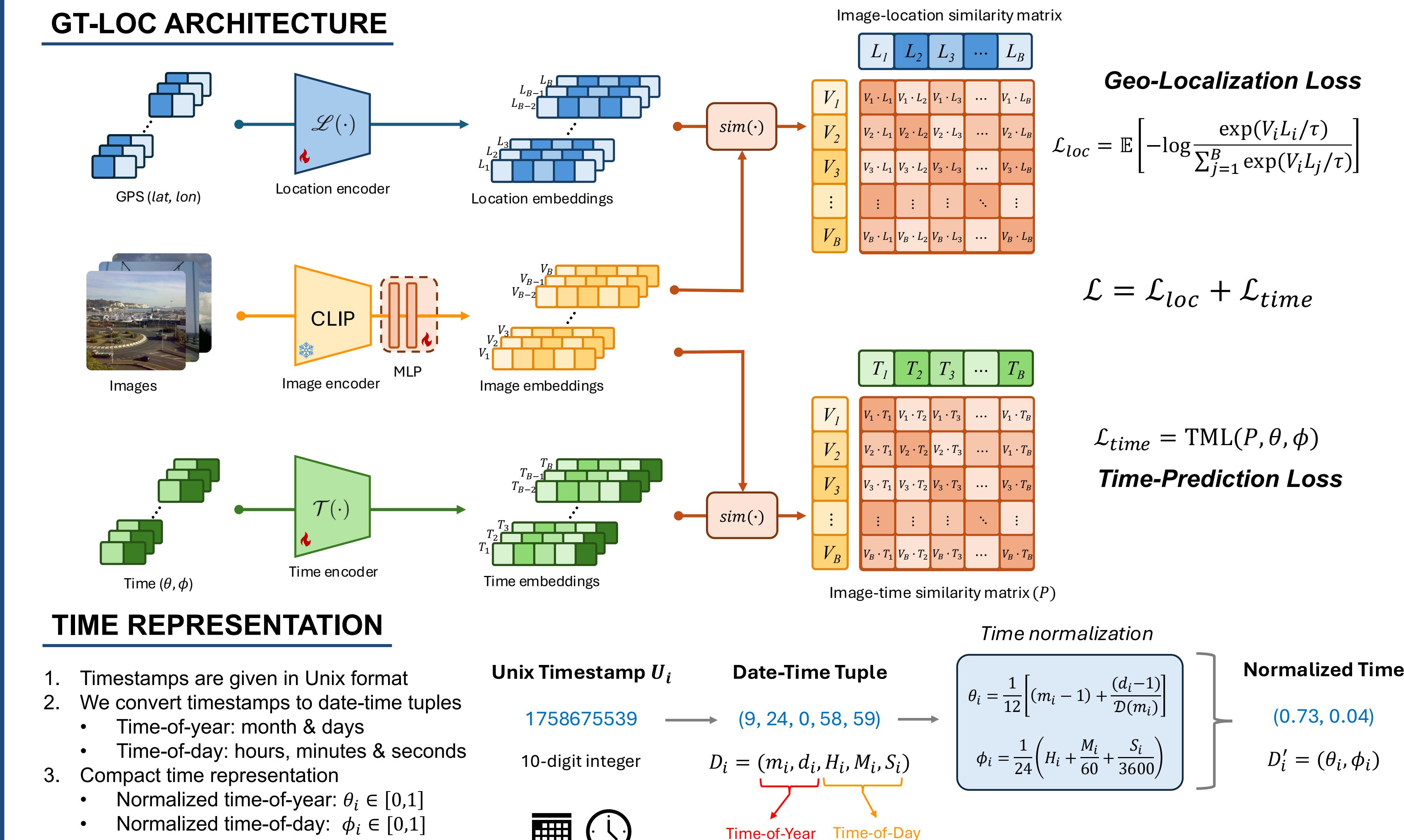


- Higher latitudes have large differences in seasons and hours of sunlight throughout the year compared to lower latitudes
- The visual appearance of time strongly depends on the geo-location

## RESULTS: TIME PREDICTION & GEO-LOCALIZATION



## GT-LOC ARCHITECTURE

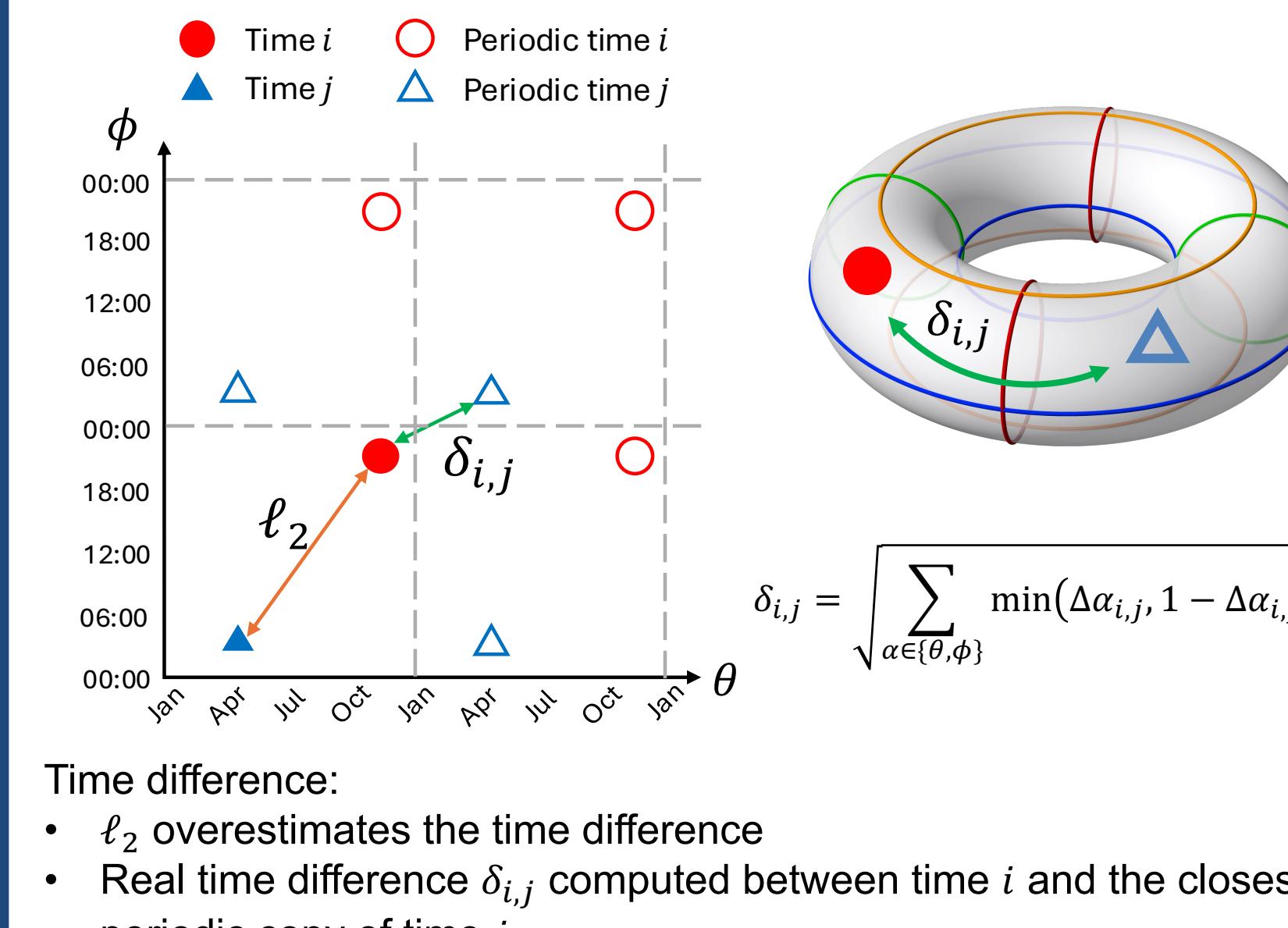


## TIME REPRESENTATION

- Timestamps are given in Unix format
- We convert timestamps to date-time tuples
  - Time-of-year: month & days
  - Time-of-day: hours, minutes & seconds
- Compact time representation
  - Normalized time-of-year:  $\theta_i \in [0, 1]$
  - Normalized time-of-day:  $\phi_i \in [0, 1]$

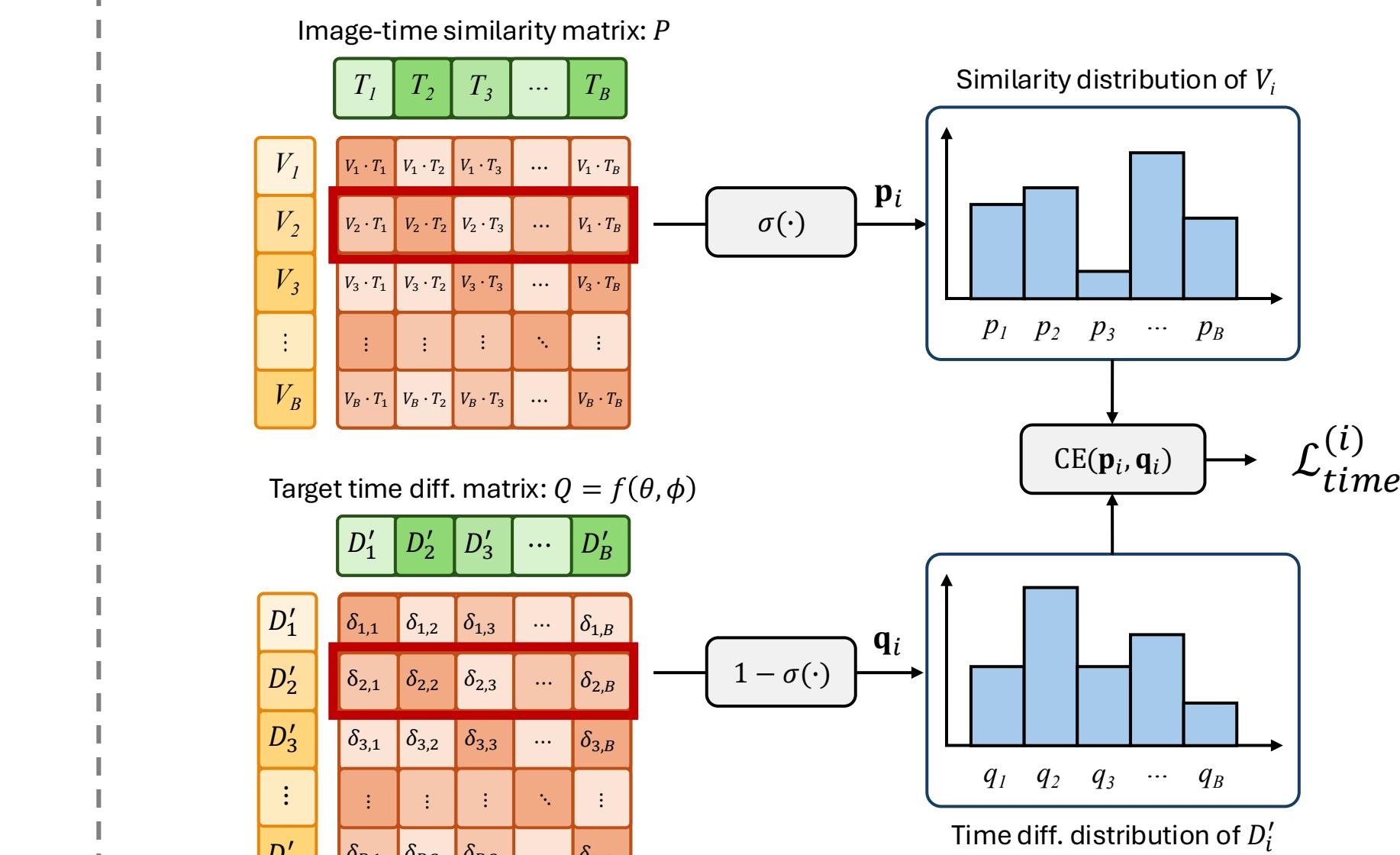
## TIME DIFFERENCE ( $\delta_{i,j}$ )

- Months and hours are periodic
- We represent time as points in month-hour plane (torus)

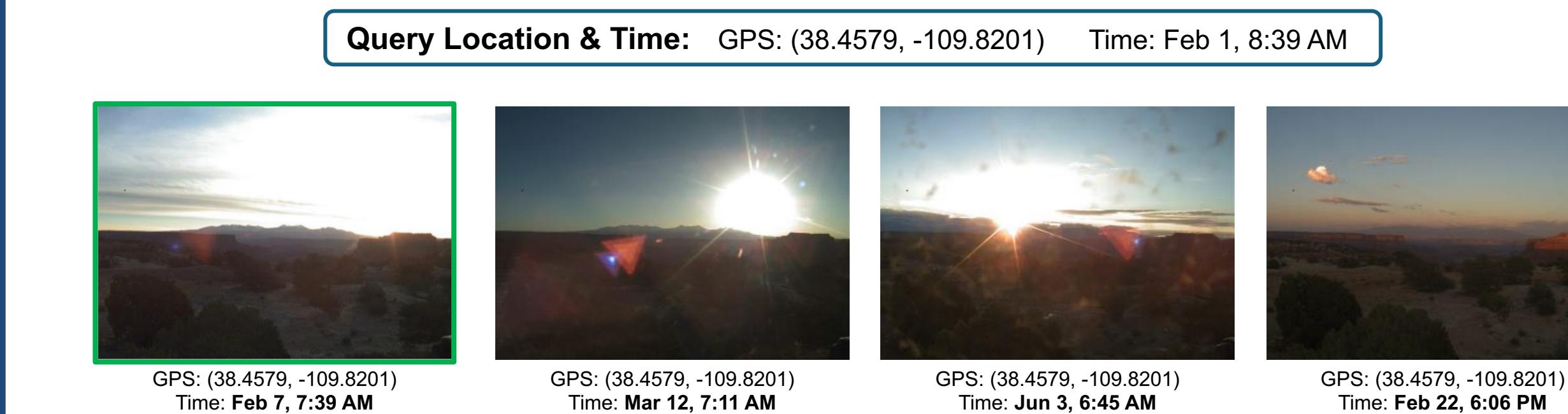


## TEMPORAL METRIC LEARNING (TML)

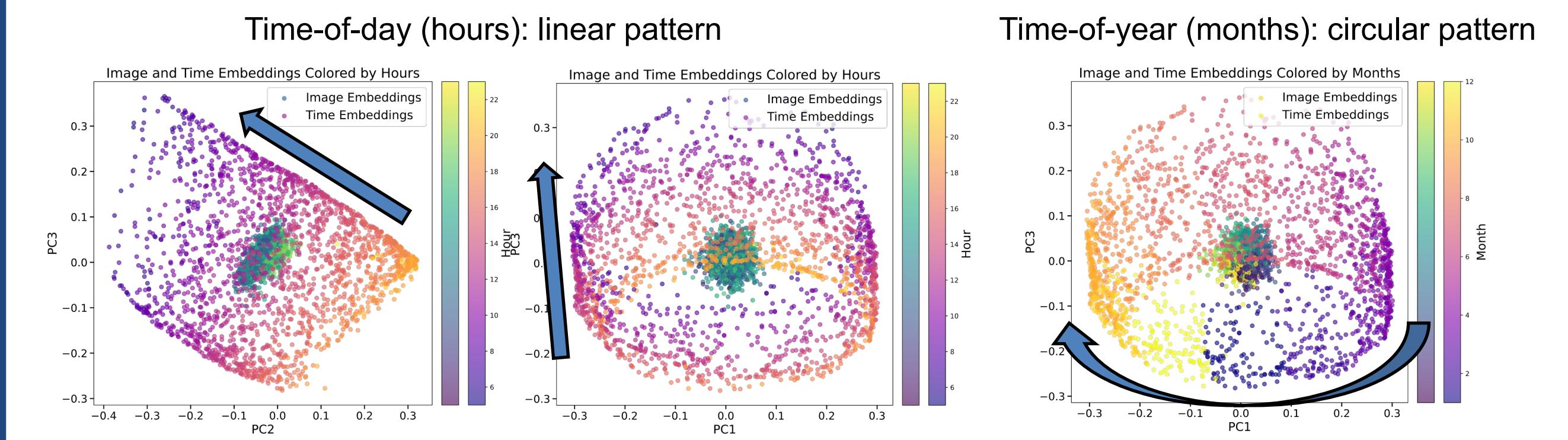
- TML: metric contrastive loss for image-time alignment
- Encourages distribution of similarities to match the distribution of time differences



## DOWNTSTREAM TASK: GEO-TEMPORAL IMAGE RETRIEVAL



## EMBEDDING SPACE ANALYSIS



## QUANTITATIVE RESULTS

SoTA performance on 3 tasks:

- Time prediction
- Geo-localization
- Geo-temporal image retrieval

Effect of joint embedding space

- GT-Loc achieves significant improvements over baselines that use GPS as input

→ Joint embedding is crucial for accurate time prediction!

Method	Im2GPS3k	GWS15k
[L] kNN, sigma=4	7.2	-
PlaNet	8.5	-
CPlaNet	10.2	-
ISNs	10.5	0.1
Translocator	11.8	0.5
GeoDecoder	12.8	0.7
GeoCLIP	14.11	0.6
PIGEOTTO	11.3	0.7
Img2Loc(LLaVA)	8.0	-
<b>GT-Loc (Ours)</b>	<b>14.41</b>	<b>0.88</b>

Geo-Temporal Image Retrieval		
Dataset: SkyFinder		
Method	R@1	R@5
Zhai et al.	0.91	7.81
Zhai et al. w/ CLIP	2.58	16.22
<b>GT-Loc (Ours)</b>	<b>6.69</b>	<b>24.58</b>
		38.54

## TAKEAWAYS

- Shared embedding space between the image, time and location modalities
  - Enables novel downstream tasks: geo-temporal image and text retrieval
- First retrieval method for time-of-capture prediction
- Novel time representation over the surface of a flat torus
- Temporal Metric Learning for image-time alignment with soft targets