

# MMLandmarks: a Cross-View Instance-Level Benchmark for Geo-Spatial Understanding

Oskar Kristoffersen<sup>1,2</sup>, Alba Reinders Sánchez<sup>1</sup>, Morten Rieger Hannemose<sup>1,2</sup>,  
Anders Bjorholm Dahl<sup>1,2</sup>, Dim P. Papadopoulos<sup>1,2</sup>



Technical University of Denmark<sup>1</sup>

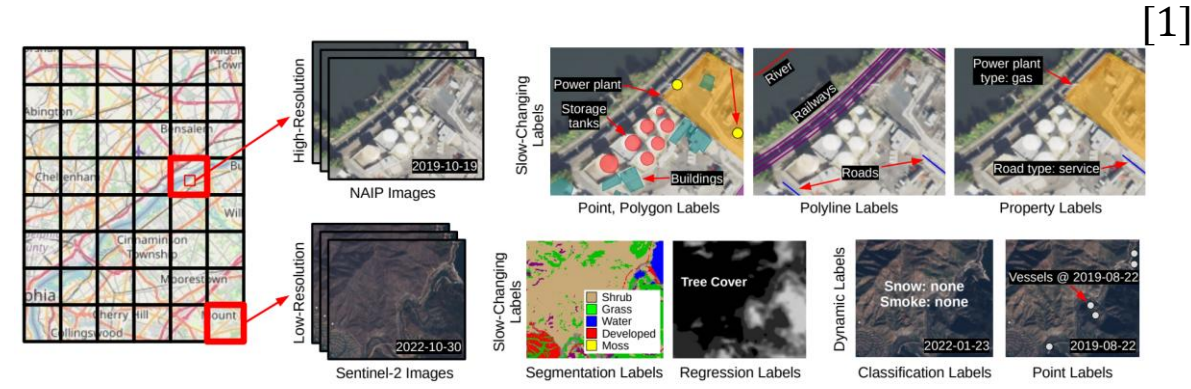


PIONEER CENTRE  
FOR ARTIFICIAL  
INTELLIGENCE

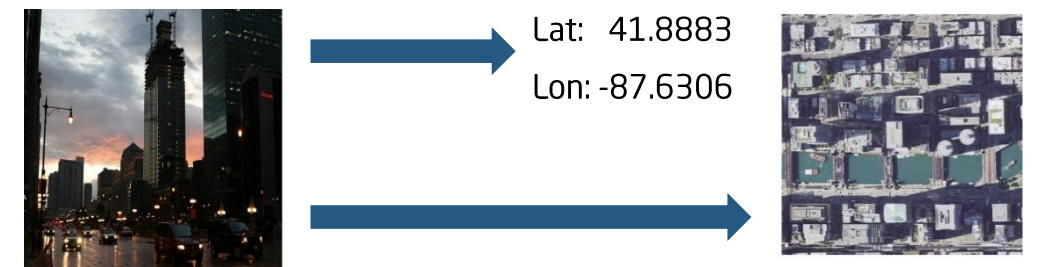
Pioneer Center for Artificial Intelligence<sup>2</sup>

# Motivation - Geospatial data

- **Multimodal data**
  - Strong embeddings
  - Foundation models for geospatial tasks



- **Aligning modalities** from the **same location**
  - Geolocalization
  - Cross-view localization

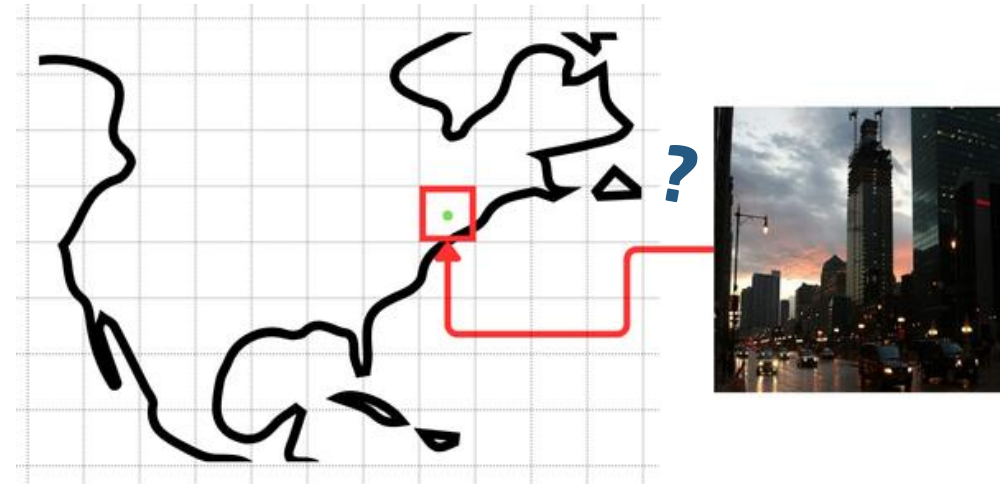
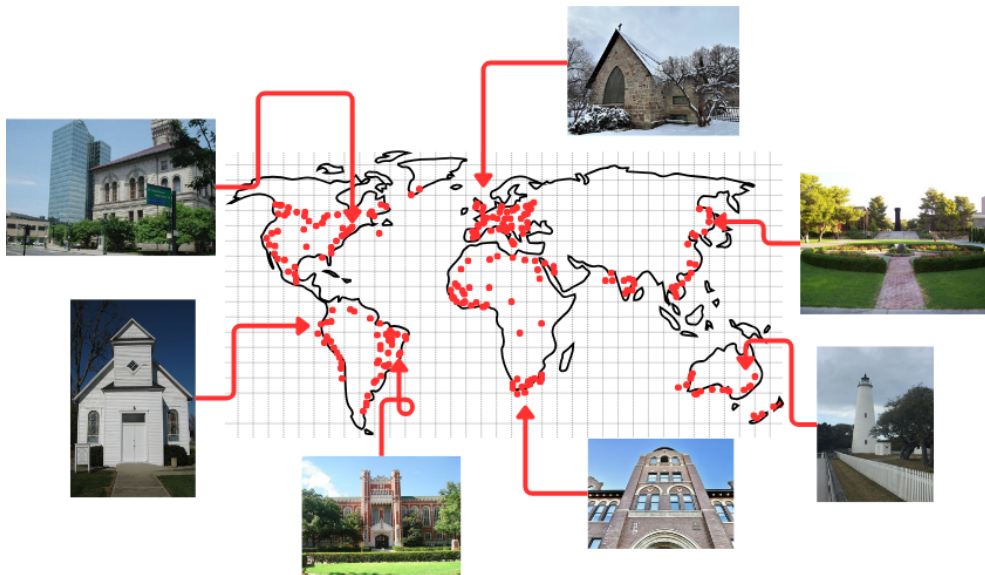


- Models solve **individual geospatial tasks** in **isolation**.
- No geospatial dataset supports **instance-level** correspondence across **multiple modalities**.

[1] Bastani et al. SatlasPretrain: A Large-Scale Dataset for Remote Sensing Image Understanding, ICCV'23

# Motivation - Geolocalization

- Cell-based classification from a grid
  - Lacks fine-grained predictions (<10km)
  - Difficult to capture enough information from a single image



# Motivation - Cross-view Localization

Current Cross-View datasets:

- 360-degree street-view images
- Have limited coverage & diversity: Streets and cities
- Build reliance on road structure



[2] Shi et al. Spatial-Aware Feature Aggregation for Cross-View Image based Geo-Localization, NeurIPS'19

[3] Ye et al. Cross-view image geo-localization with Panorama-BEV Co-Retrieval Network, ECCV'25

# Combining all modalities into one dataset

### Space Needle

The Space Needle is an observation tower in Seattle, Washington...

lat : 47.6205  
lon: -122.3493

Wikimedia Commons: [Image 1] ... [Image 2] ... [Image 3]

GEE: [Image 1] (2007) ... [Image 2] (2012) ... [Image 3] (2016)

### Liberty Island

Liberty Island is a federally owned island in Upper New York Bay...

lat : 40.6898  
lon: -74.0453

Wikimedia Commons: [Image 1] ... [Image 2] ... [Image 3]

GEE: [Image 1] (2008) ... [Image 2] (2011) ... [Image 3] (2016)

### Grand Canyon Skywalk

The Grand Canyon Skywalk is a horseshoe-shaped cantilever...

lat : 36.0120  
lon: -113.8111

Wikimedia Commons: [Image 1] ... [Image 2] ... [Image 3]

GEE: [Image 1] (2006) ... [Image 2] (2009) ... [Image 3] (2013)

### Bristol Station

Bristol station (locally known as Union Station and Bristol Train...

lat : 36.5956  
lon: -82.1799

Wikimedia Commons: [Image 1] ... [Image 2] ... [Image 3]

GEE: [Image 1] (2003) ... [Image 2] (2007) ... [Image 3] (2017)

### Old Citrus County

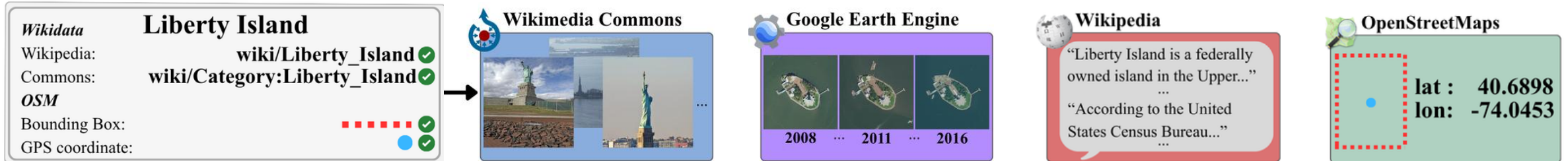
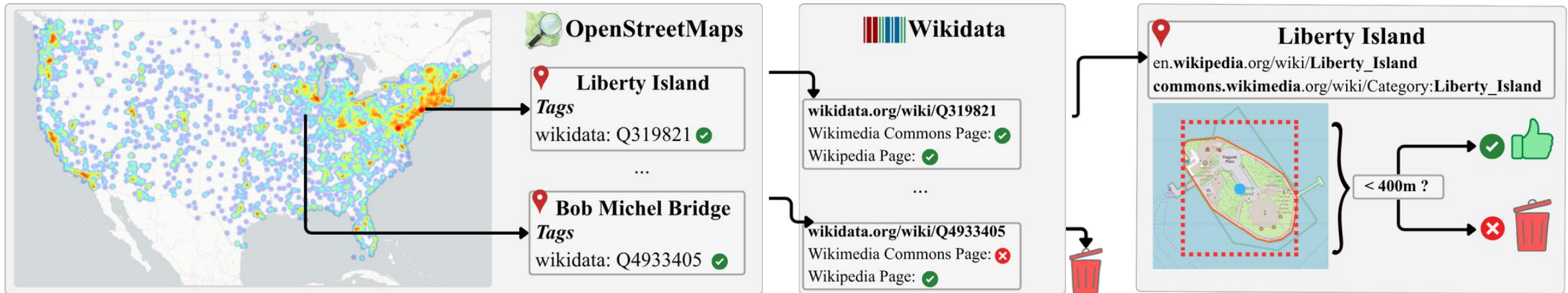
The Old Citrus County Courthouse (constructed in 1912) is a historic...

lat : 28.8360  
lon: -82.3302

Wikimedia Commons: [Image 1] ... [Image 2] ... [Image 3]

GEE: [Image 1] (2005) ... [Image 2] (2008) ... [Image 3] (2018)

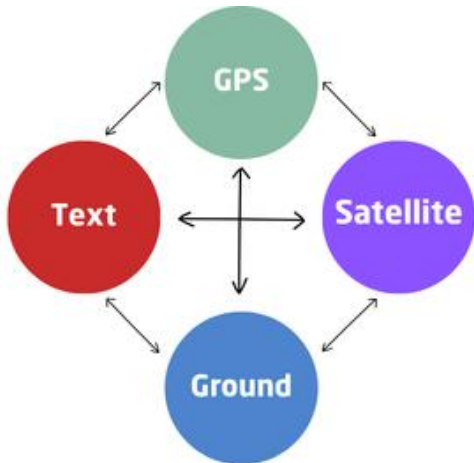
# Our pipeline for data collection





# Final dataset overview

- 18557 Landmarks across the United States of America



329k Ground Images



197k Satellite Images



18,557 Coordinates

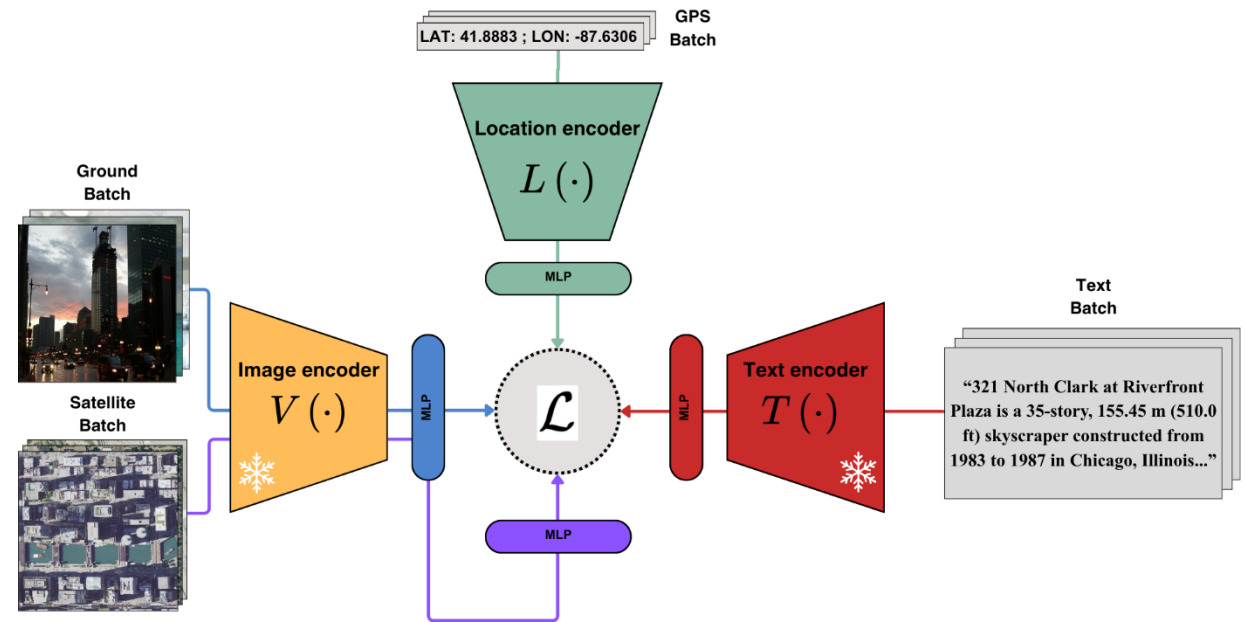


18,557 Text Information

# A simple baseline

- Frozen CLIP Image and Text encoders.
- Initialized GeoCLIP [1] Location encoder
- Trainable MLP Projection heads
- Complete Contrastive Loss

$$\mathcal{L} = \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{j=1, j \neq i}^K \mathcal{L}_{i,j}$$



# Some Results

|                  | model          | Satellite → Ground |             |             | Ground → Satellite |             |             |
|------------------|----------------|--------------------|-------------|-------------|--------------------|-------------|-------------|
|                  |                | medR ↓             | R@1 ↑       | R@5 ↑       | medR ↓             | R@1 ↑       | R@5 ↑       |
| <i>Zero-Shot</i> | DiNOv2         | 28342              | 1.9         | 3.0         | 9335               | 3.3         | 8.8         |
|                  | DiNOv3         | 14623              | 3.7         | 7.4         | 2241               | 6.1         | 15.5        |
|                  | SigLIP         | 1546               | 7.2         | 15.25       | 679                | 10.7        | 23.8        |
|                  | OAI-CLIP       | <u>519</u>         | <u>10.4</u> | <u>21.2</u> | 620                | 11.9        | 26.4        |
|                  | SigLIP2        | 682                | 8.6         | 18.6        | <u>140</u>         | <u>14.6</u> | <u>31.6</u> |
|                  | Uni-1652       | 62606              | 0.4         | 0.8         | 39961              | 0.2         | 0.9         |
|                  | TransGeo-90°   | 40973              | 0.7         | 2.0         | 13425              | 0.4         | 1.7         |
|                  | Sample4Geo-UNI | 34988              | 3.0         | 6.2         | 40056              | 0.85        | 2.8         |
|                  | <b>MMCLIP</b>  | <b>23</b>          | <b>30.4</b> | <b>52.4</b> | <b>48</b>          | <b>20.5</b> | <b>46.5</b> |

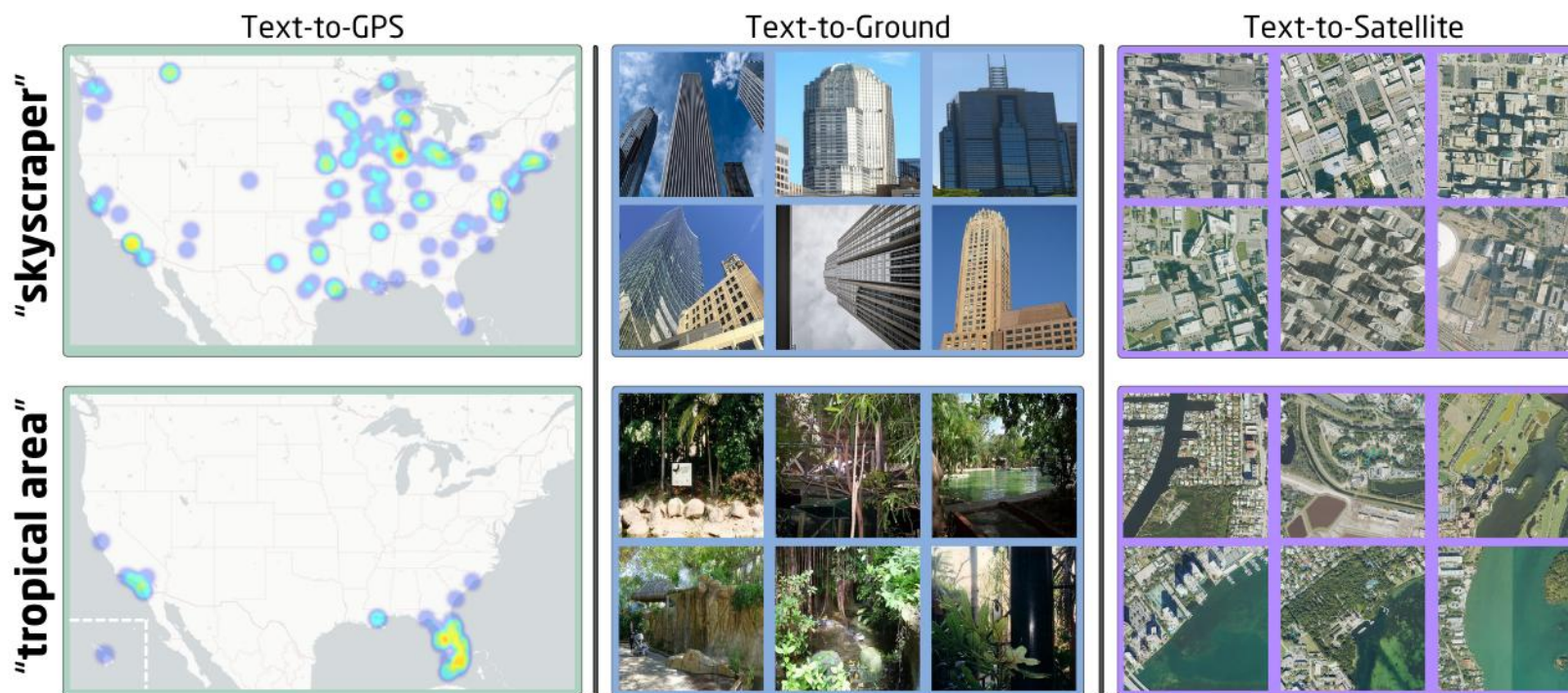
### Cross-View Localization

|                  | Method        | Distance (% @ km) ↑ |              |              |                  |             |             |
|------------------|---------------|---------------------|--------------|--------------|------------------|-------------|-------------|
|                  |               | Street<br>1 km      |              |              | Region<br>200 km |             |             |
| <i>Zero-Shot</i> |               | Ground → GPS        |              |              | Satellite → GPS  |             |             |
|                  | GeoReasoner   | 0.01                | 8.03         | 80.28        | 0.6              | 6.2         | 82.9        |
|                  | StreetCLIP    | 0.53                | 48.89        | 88.48        | 0.6              | 49.6        | 93.9        |
|                  | osv5m         | 1.28                | 19.19        | 57.25        | 0.7              | 21.2        | 70.8        |
|                  | G3            | 12.86               | 45.52        | <b>91.83</b> | 8.8              | 47.9        | 97.0        |
|                  | GeoCLIP       | <b>21.37</b>        | 48.57        | 91.5         | <u>12.3</u>      | <u>48.8</u> | <u>97.4</u> |
|                  | <b>MMCLIP</b> | 16.83               | <b>51.78</b> | <u>91.52</u> | <b>36.9</b>      | <b>81.1</b> | <b>99.7</b> |

### Geo-localization

- Comparison with off the shelf foundation and specialized models
- Goal: Illustrate flexibility, not superiority
- MMCLIP, is versatile when trained on Mmlandmarks.

# Text-to-X modality alignment





# Final Thoughts - MMLandmarks

- First large-scale, instance-level, multimodal dataset for geospatial understanding
- 18557 landmarks across the United States
- Landmark-level one to one correspondence: Ground images, Aerial images, Text and GPS
- Open Domain - permissive licenses

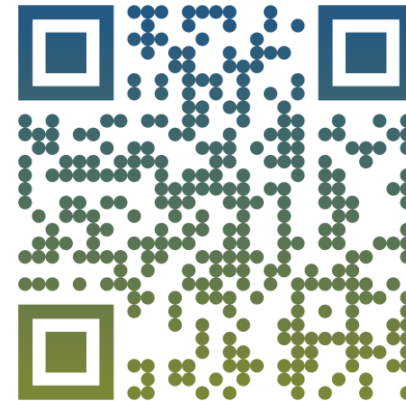
Dataset



Paper



Codebase



**Project page - scan me!**  
[mmlandmarks.compute.dtu.dk](https://mmlandmarks.compute.dtu.dk)

DTU

