ETH zürich [1]

Meta [2]

REALITY LABS

# OrienterNet

## Visual Localization in 2D Public Maps with Neural Matching

Paul-Edouard Sarlin[1]   Daniel DeTone[2]   Tsun-Yi Yang[2]   Armen Avetisyan[2]

Julian Straub[2]   Tomasz Malisiewicz[2]   Samuel Rota Bulo[2]

Richard Newcombe[2]   Peter Kontschieder[2]   Vasileios Balntas[2]
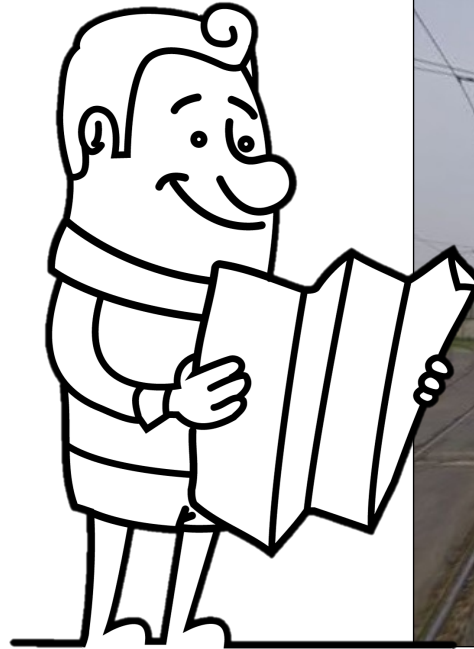
CVPR 2023        psarlin.com/orienternet        Poster THU-PM-098

# Humans use simple 2D maps
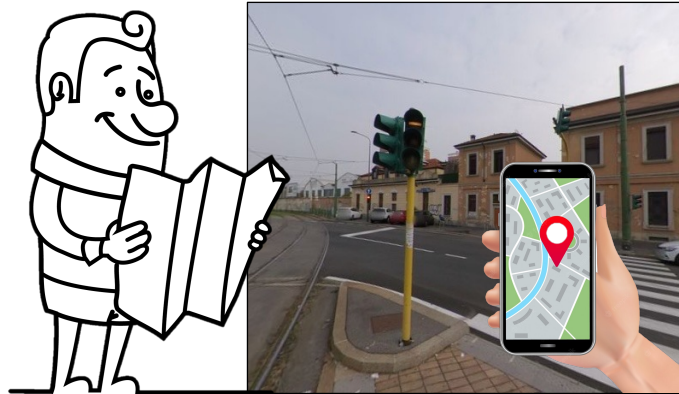
where am I?

input image
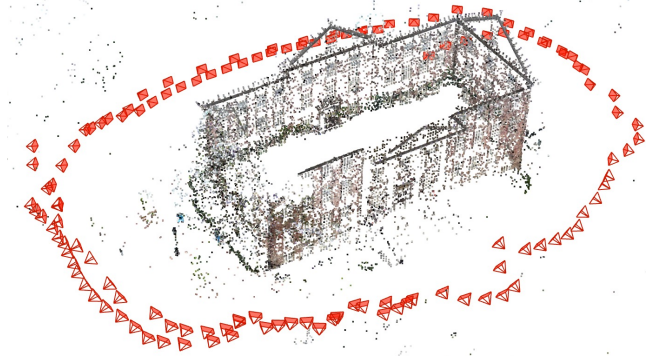
*inputs*

input image

gravity
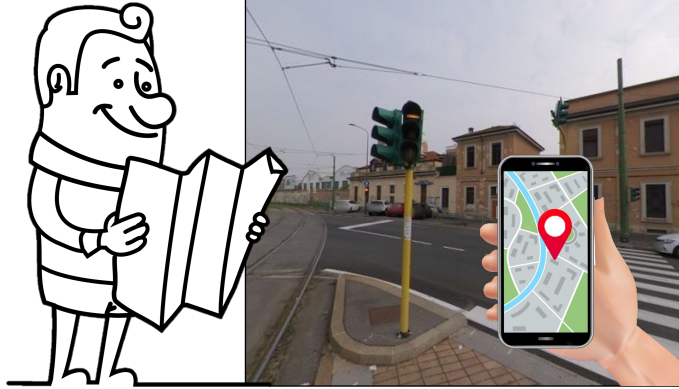
Humans use 2D maps

Existing algorithms: 3D point clouds
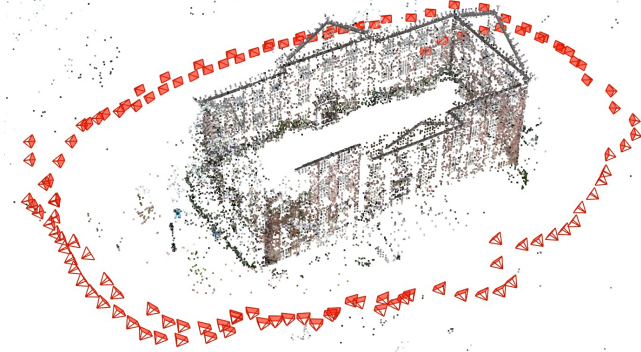
*inputs*

input image

gravity
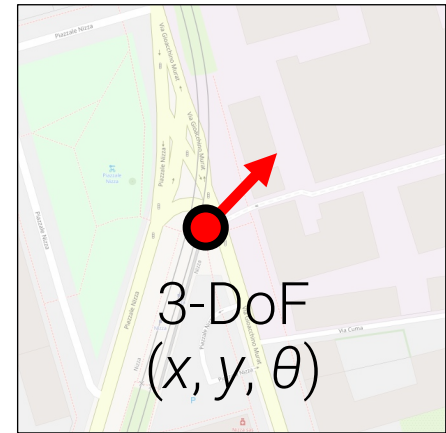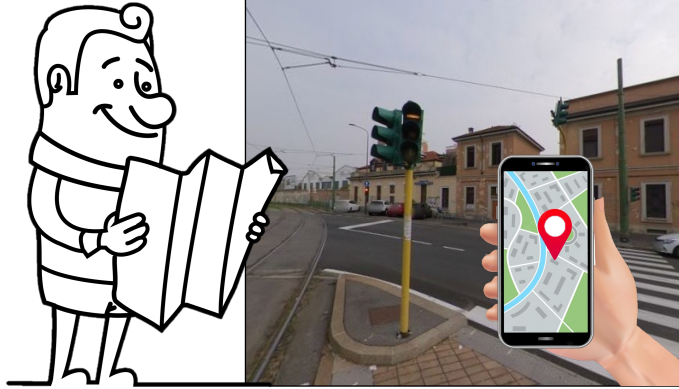
Humans use 2D maps

Existing algorithms: 3D point clouds

GPS

OpenStreetMap

**OrienterNet**

3-DoF $(x, y, \theta)$
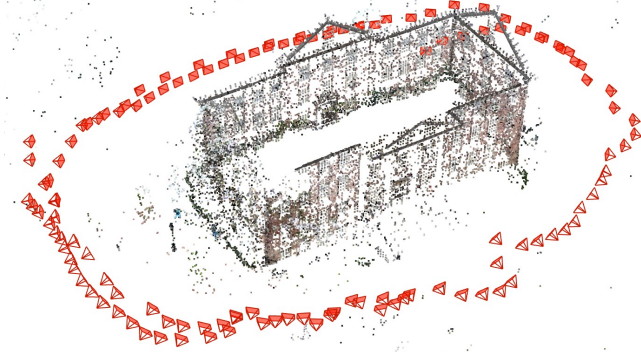
camera pose

*inputs*

input image

gravity

Humans use 2D maps

Existing algorithms: 3D point clouds

**OrienterNet**
neural map matching

GPS

OpenStreetMap

3-DoF
$(x, y, \theta)$

camera pose

*inputs*

input image

gravity

**Humans use 2D maps**

**Existing algorithms: 3D point clouds**
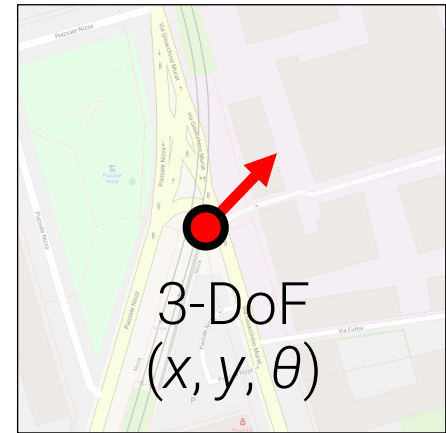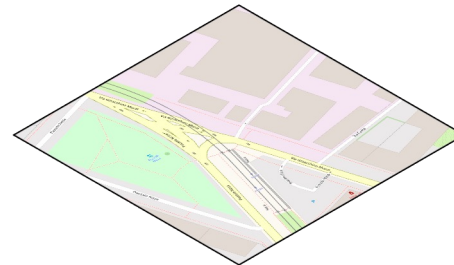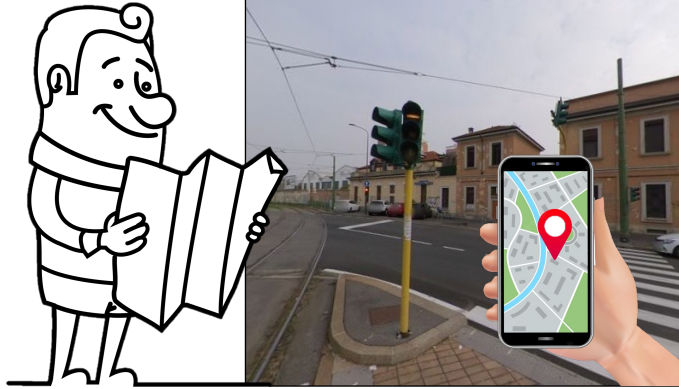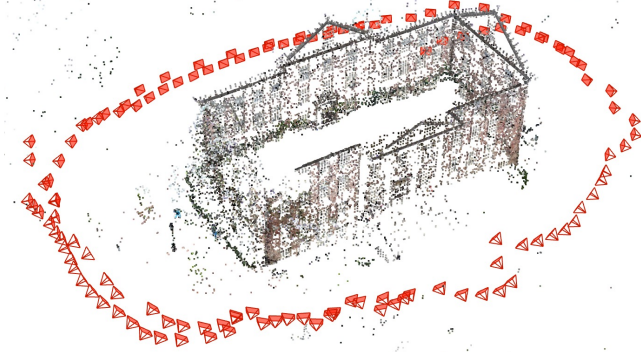
GPS

OpenStreetMap

**OrienterNet**
neural map matching

3-DoF
$(x, y, \theta)$

camera pose

# inputs

input image

gravity

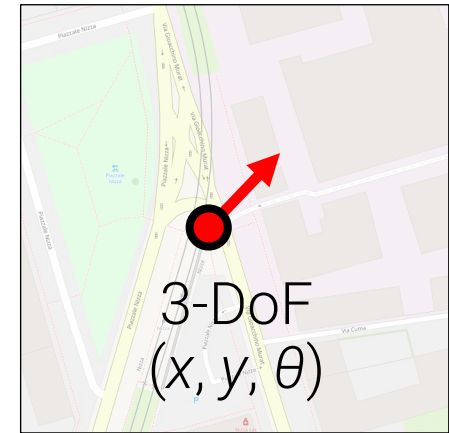## Humans use 2D maps

## Existing algorithms: 3D point clouds

## OrienterNet
neural map matching

OpenStreetMap

3-DoF
$(x, y, \theta)$
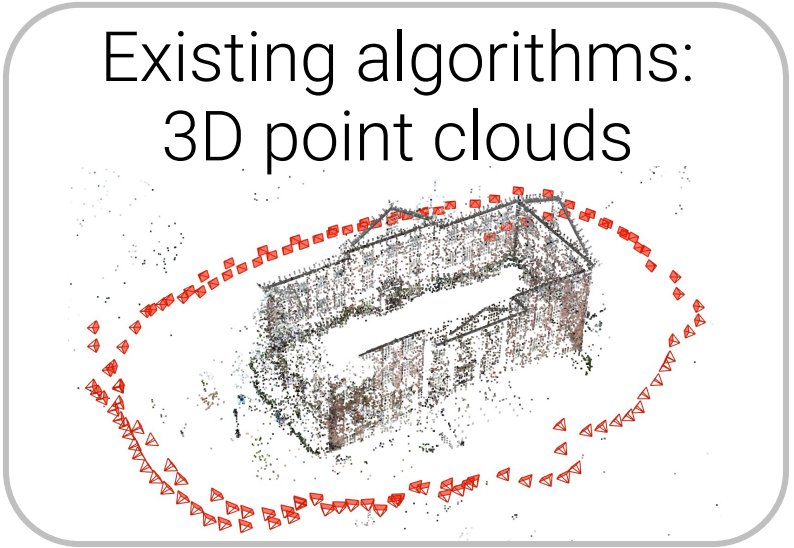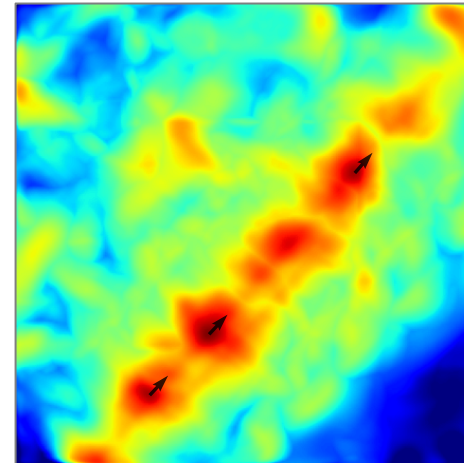
camera pose

# Zero-shot generalization



Aria

Mapillary

KITTI

# Zero-shot generalization

Aria 

Mapillary 

KITTI



6

# Positioning

Recover the 6-DoF pose of the device
- 3D translation + rotation
- global reference frame



pose

# GPS+compass is not enough

- Low accuracy
- Only 3 DoF
- Commonly unreliable: urban canyon, metal structures



Actual Position



Google Maps

R, t ?

6-DoF Localization

R, t ?

6-DoF Localization

R, t ?

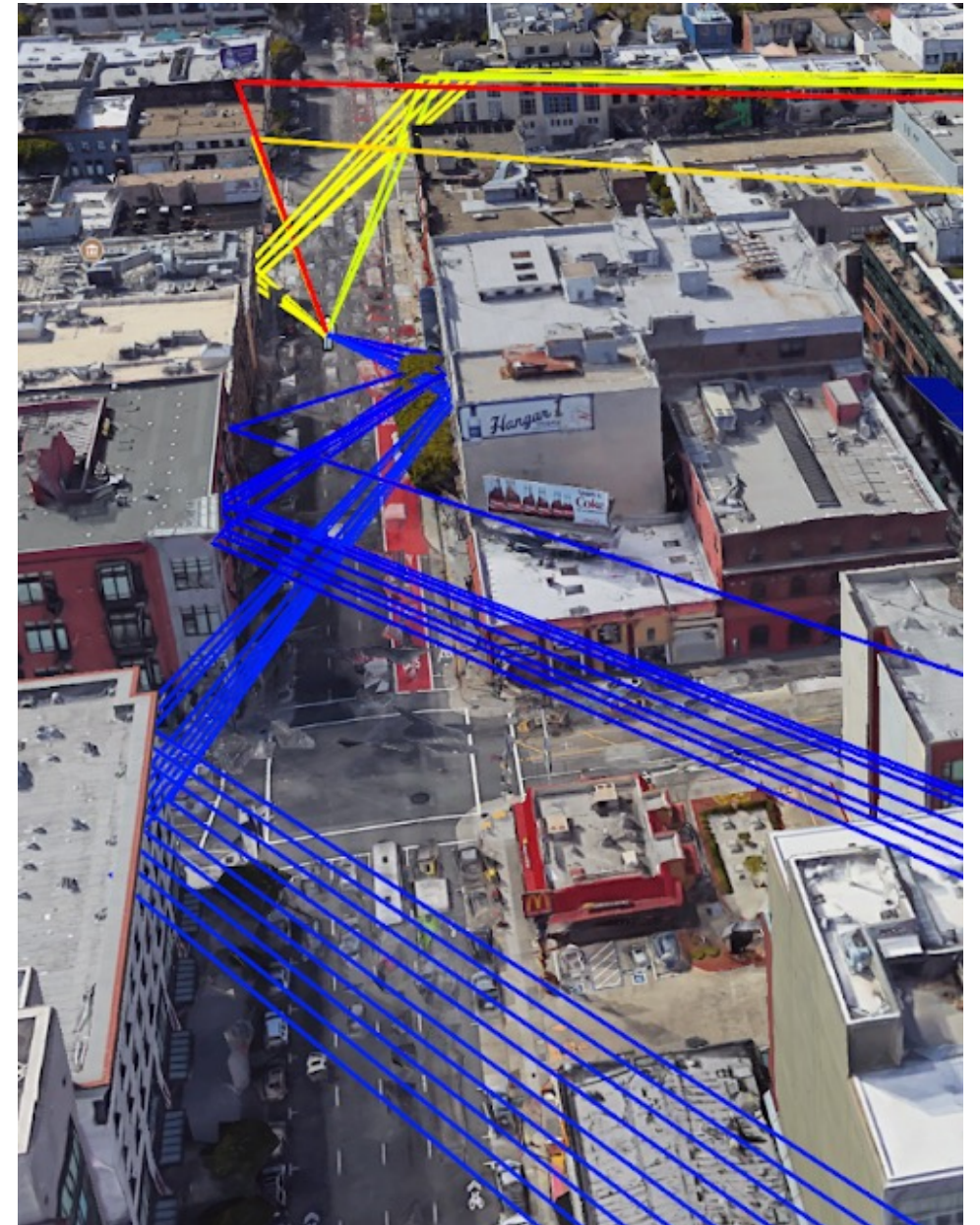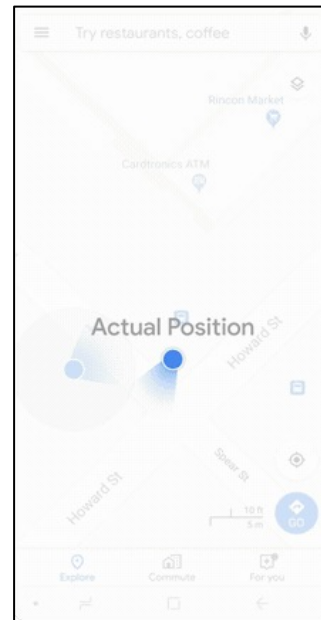6-DoF Localization

Structure-from-Motion

R, t ?

6-DoF Localization

10

R, t ?

6-DoF Localization

# Limitations of 3D maps

Build
& update

**Mapping fleet**
**Frequent updates**

Google StreetView

Storage

Privacy

Mapillary

# Limitations of 3D maps

🔄 Build
& update

🗄 **Storage**                Very large

🛡 Privacy

# Limitations of 3D maps



Build & update

Storage

**Privacy**

keypoints    descriptors    reconstruction

inversion

Mihai Dusmanu

Risk of inversion

# Limitations of 3D maps

 Build & update — Mapping fleet / Frequent updates

Compression & Quantization

 Storage — Very large

 Privacy — Risk of inversion

content-concealing descriptors

Privacy-preserving descriptors

# Semantic 2D maps

Planimetric

OpenStreetMap

# Rich semantics



building boundary

ticket machine

pharmacy

bench

trash can

post box

bike parking

restrooms

tree

pedestrian path

# Rich semantics



Google StreetView

building boundary

ticket machine

pharmacy

bench

trash can

bike parking

post box

restrooms

tree

pedestrian path

16

|  | 3D maps | 2D maps |
|---|---|---|
| Build & update | Mapping fleet Frequent updates | Public No appearance updates |
| Storage | Very large | Compact Transfer on-device |
| Privacy | Risk of inversion | No private info |

# Simplifying assumptions



- Known gravity direction

- Unnecessary vertical position

3-DoF pose
$(x, y, \theta)$

# Problem setup



image
+ gravity

128m x 128m



OpenStreetMap

The
OrienterNet
architecture

1. Bird's-Eye
View inference

image

g

2. Map encoding

GPS

OpenStreetMap

3-DoF pose
$(x, y, \theta)$

20

# The OrienterNet architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

**2. Map encoding**

OpenStreetMap

GPS

3-DoF pose
$(x, y, \theta)$

# The OrienterNet architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

**2. Map encoding**

GPS

OpenStreetMap

3-DoF pose
$(x, y, \theta)$

# The OrienterNet architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

**2. Map encoding**

OpenStreetMap

GPS

Rasterize

raster map

3-DoF pose $(x, y, \theta)$

# The OrienterNet architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

**2. Map encoding**

OpenStreetMap

GPS

Rasterize

raster map

CNN

neural map

**3-DoF pose** $(x, y, \theta)$

# The OrienterNet architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

**2. Map encoding**

OpenStreetMap

GPS

Rasterize

raster map

CNN

neural map

**3-DoF pose** $(x, y, \theta)$

rotate + correlate

**3. BEV-map matching**

The
OrienterNet
architecture



1. Bird's-Eye
View inference

image

g

CNN

occupancy
volume

pool + CNN

neural BEV

confidence

2. Map encoding

GPS

OpenStreetMap

Rasterize

raster map

CNN

neural map

3-DoF pose
$(x, y, \theta)$

$\theta$

$x$

$y$

pose likelihood

rotate + correlate

3. BEV-map
matching

The
OrienterNet
architecture



**1. Bird's-Eye View inference**

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

**2. Map encoding**

GPS

OpenStreetMap

Rasterize

raster map

CNN

neural map

**3-DoF pose** $(x, y, \theta)$

argmax

$\theta$  $x$  $y$

pose likelihood

rotate + correlate

**3. BEV-map matching**

20

The OrienterNet architecture

1. Bird's-Eye View inference

image

g

CNN

occupancy volume

pool + CNN

neural BEV

confidence

2. Map encoding

OpenStreetMap

GPS

Rasterize

raster map

CNN

neural map

3-DoF pose $(x, y, \theta)$

argmax

$\theta$

$x$

$y$

pose likelihood

rotate + correlate

3. BEV-map matching

Maximize likelihood of GT pose

# 1. Bird's Eye View inference

# 1. Bird's Eye View inference



gravity-aligned

# 1. Bird's Eye View inference



column

BEV

depth planes $d$

polar ray

# 1. Bird's Eye View inference



pixel scales $\boldsymbol{\sigma}$

column

camera intrinsics

**BEV**

depth planes $d$

polar ray

# 1. Bird's Eye View inference



column

pixel scales $\sigma$

BEV

polar ray

distant

expected depth
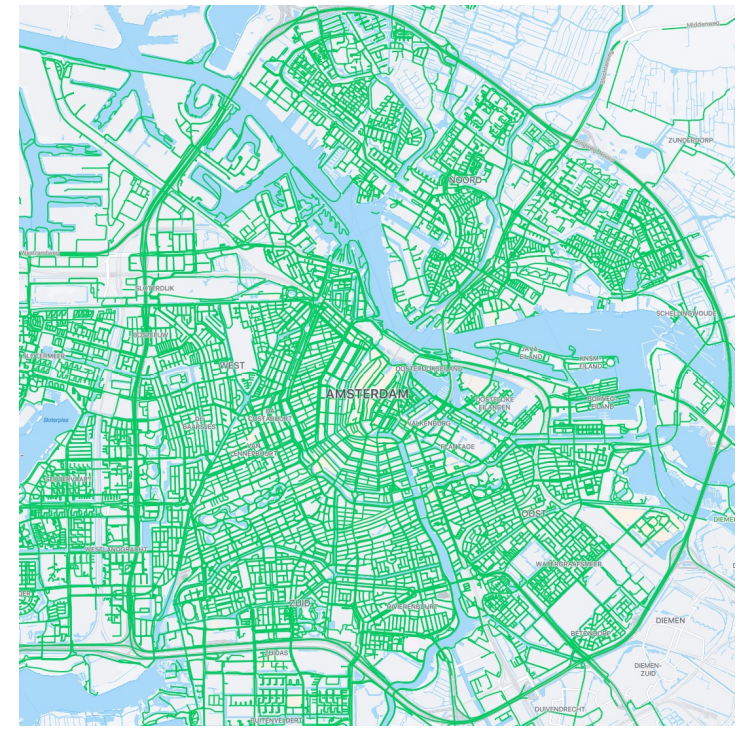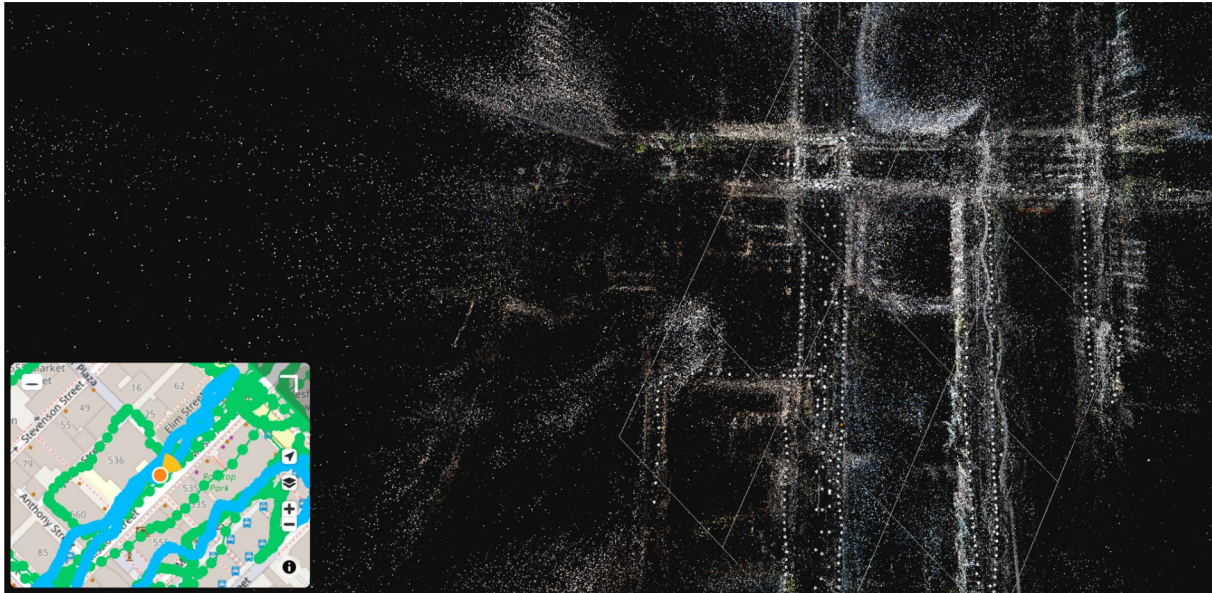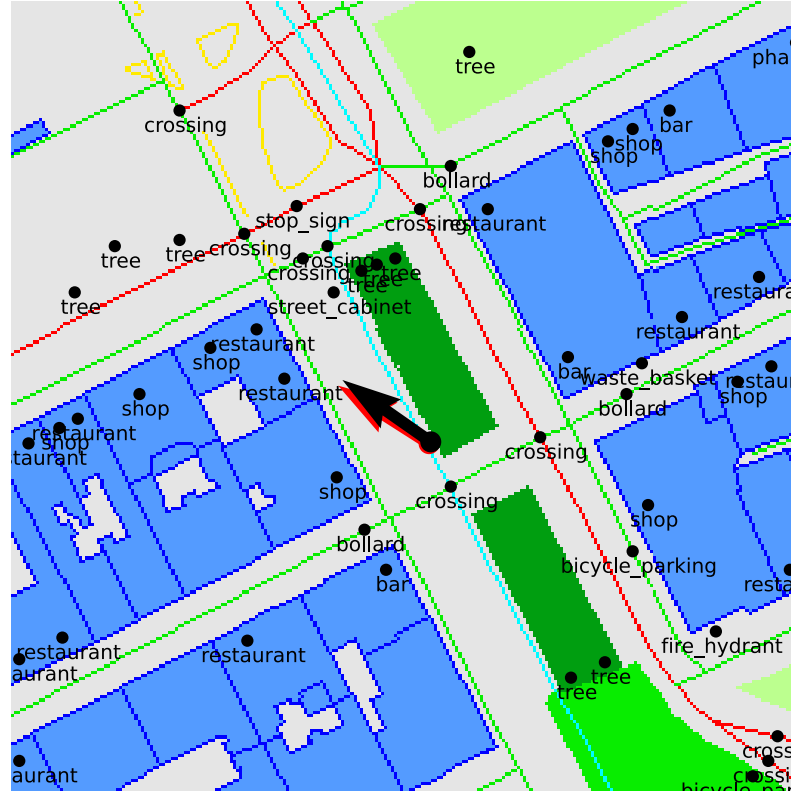
near

# Training a single strong model



- Publicly-available data from Mapillary

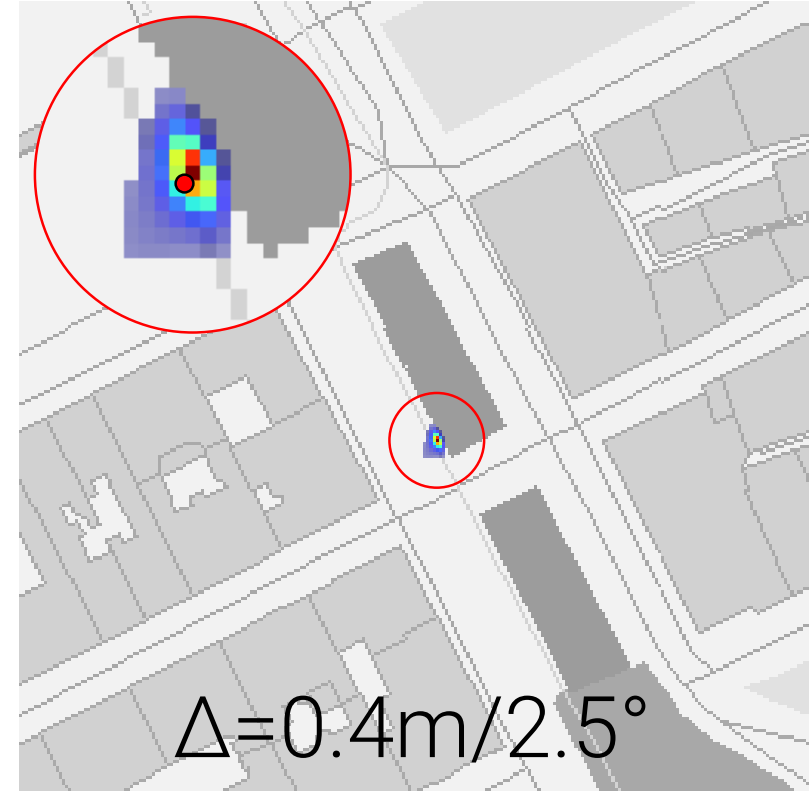- 760k images from 12 cities across Europe & US

- Hand-held, car, bike

# input image

# raster map

# likelihood



ground truth

prediction

Δ=0.4m/2.5°

building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence

# input image

# raster map

# likelihood


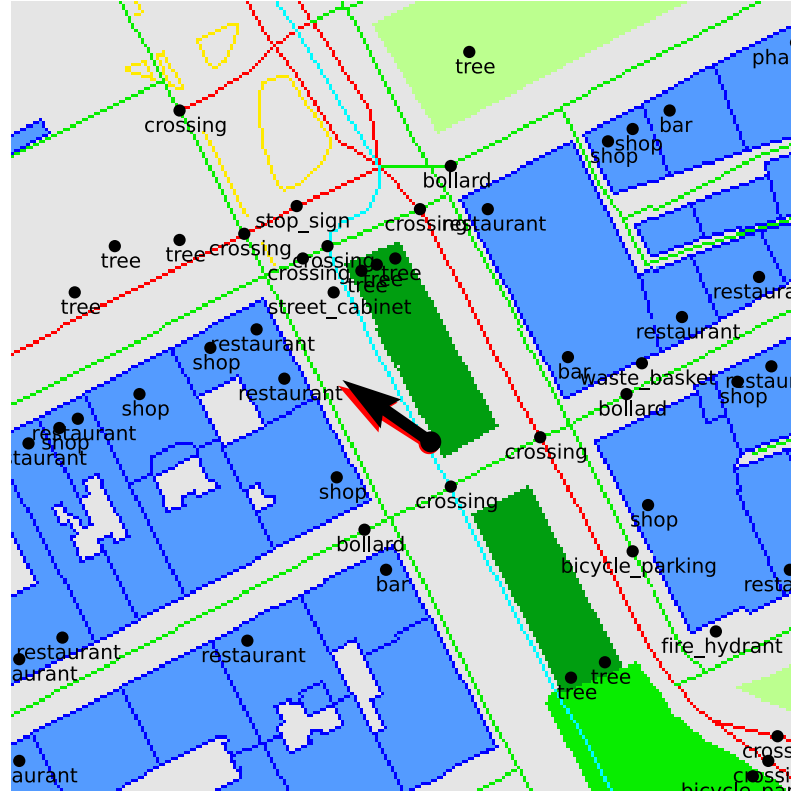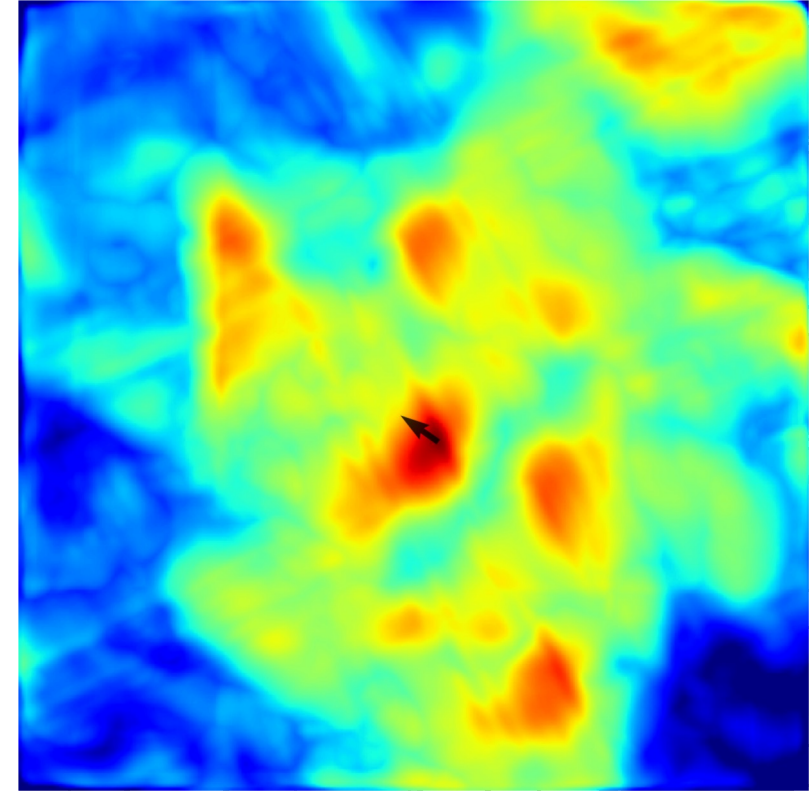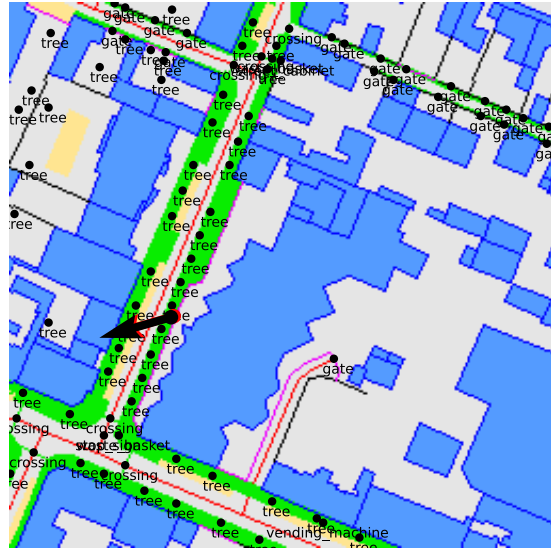
ground truth  prediction

building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence
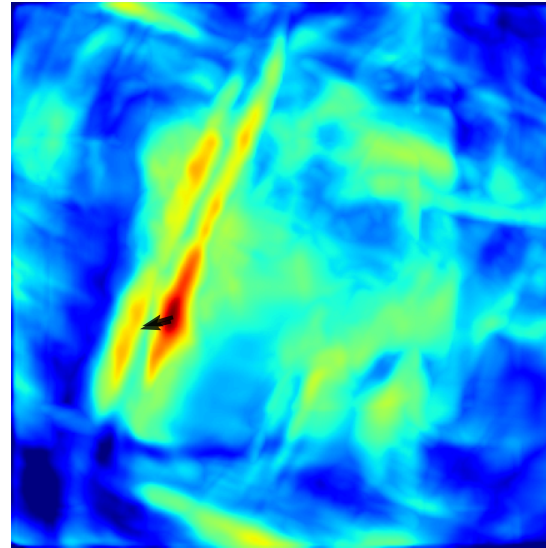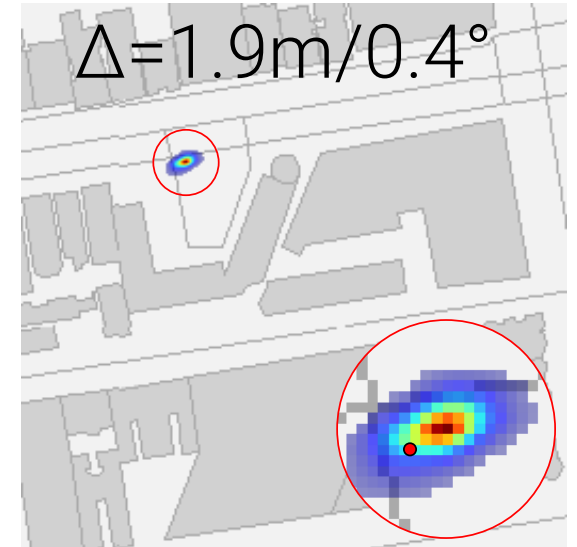
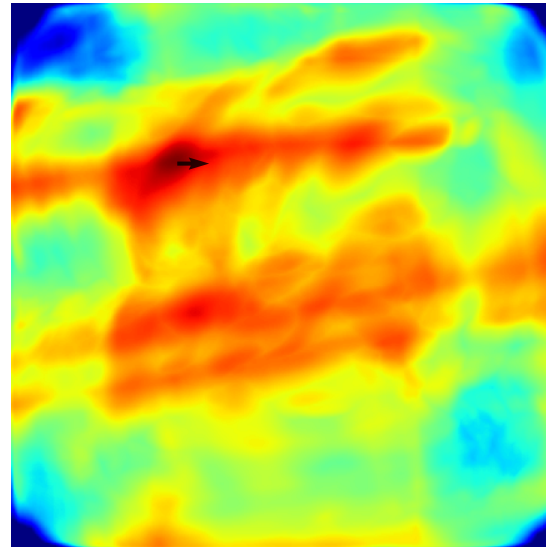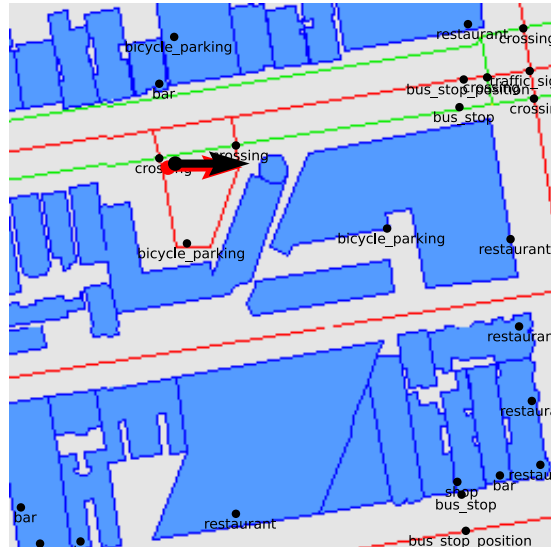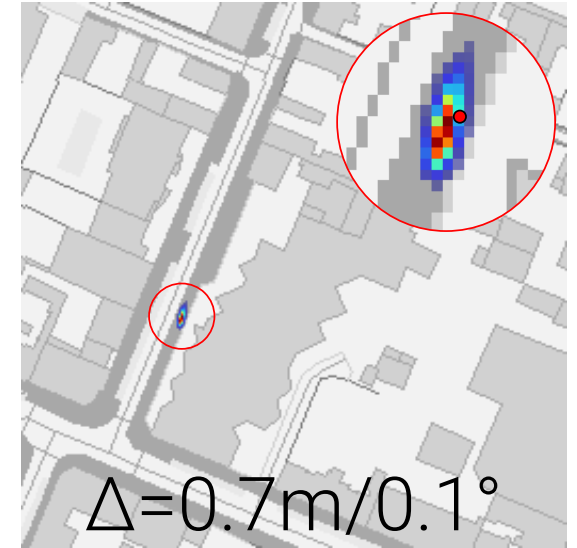| input image | raster map | log-likelihood | likelihood |
|---|---|---|---|

Δ=0.7m/0.1°

Δ=1.9m/0.4°

building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence

# Driving data – KITTI



Δ=0.7m/1.7° lateral=0.6m    long.=0.3m

Δ=0.9m/1.7° lateral=0.3m    long.=0.9m

building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence

# Driving data – KITTI



Δ=0.7m/1.7° lateral=0.6m    long.=0.3m

Δ=0.9m/1.7° lateral=0.3m    long.=0.9m
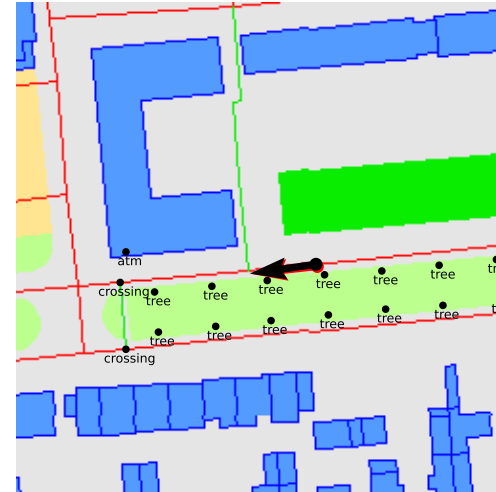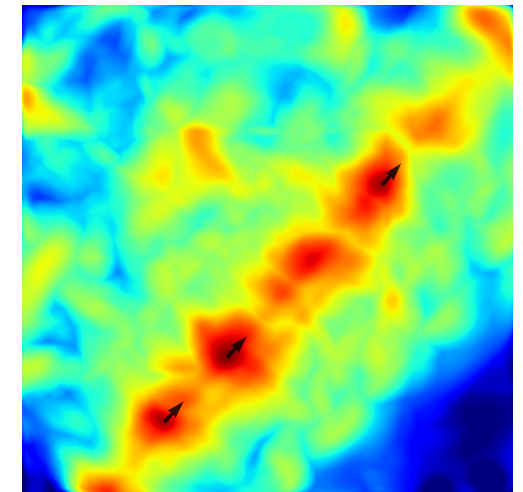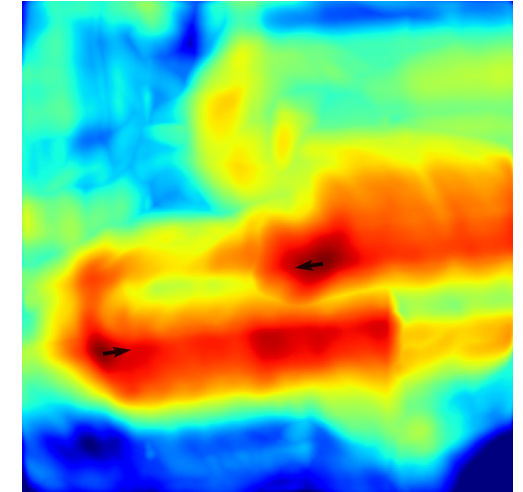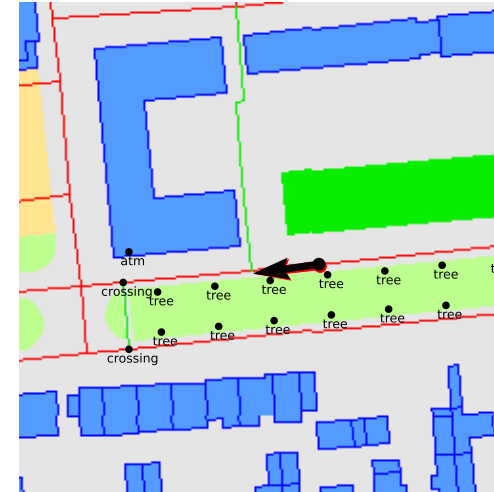
building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence

# AR data – Aria glasses



Seattle

Detroit

GPS

Δ=0.4m/5.7°

Δ=1.2m/0.2°

building ●area and ●outline, ●road, ●footway, ●cycleway, ●grass, ●park, ●playground, ●parking, ●fence
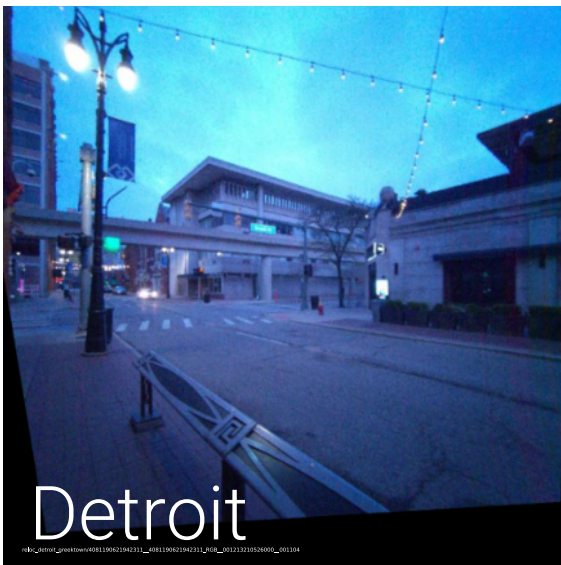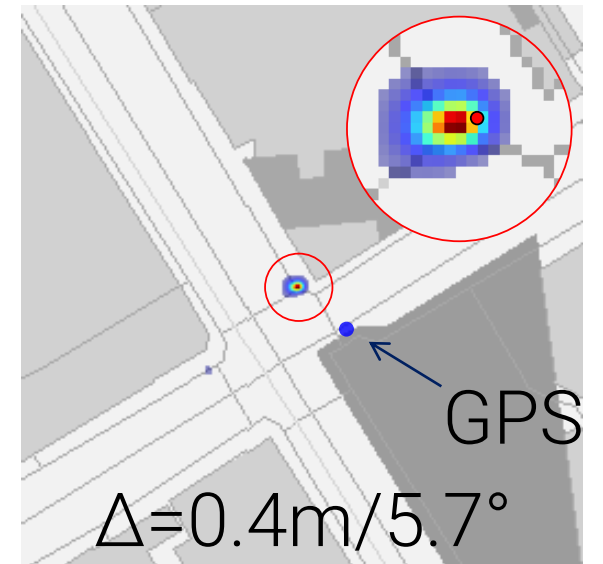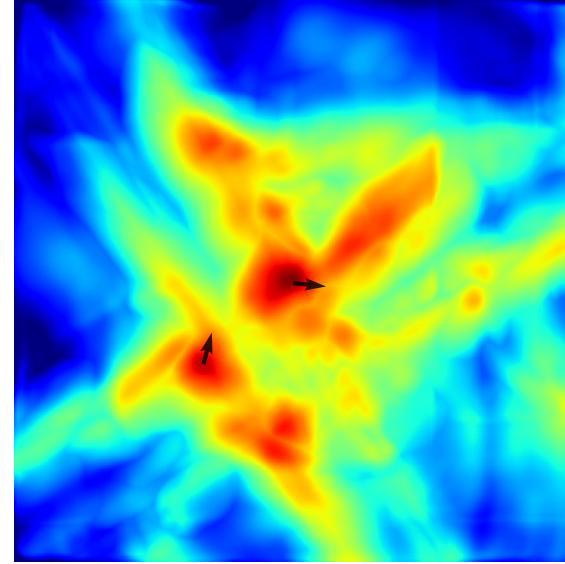
# Sequence localization

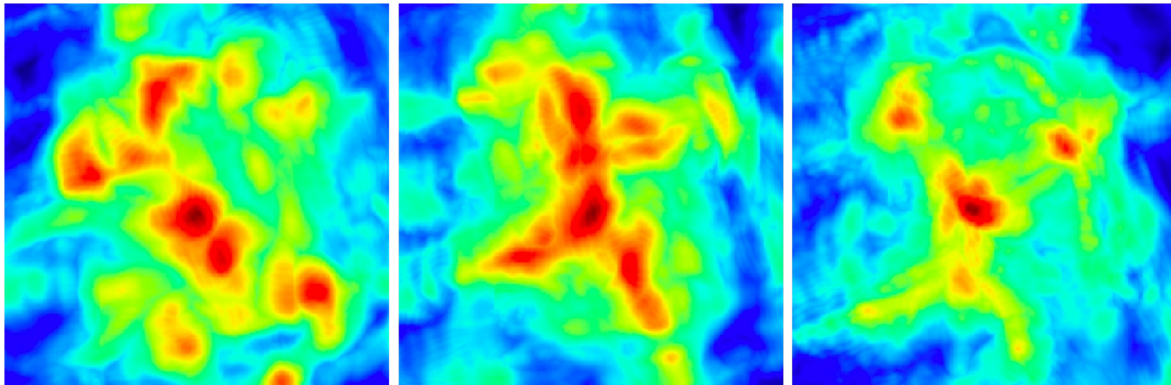**Fuse successive predictions** assuming known relative poses

$$P(\boldsymbol{\xi}_i | \{\mathbf{I}_j\}, \mathrm{map}) = \prod_k P(\boldsymbol{\xi}_i \oplus \hat{\boldsymbol{\xi}}_{ij} | \mathbf{I}_j, \mathrm{map})$$
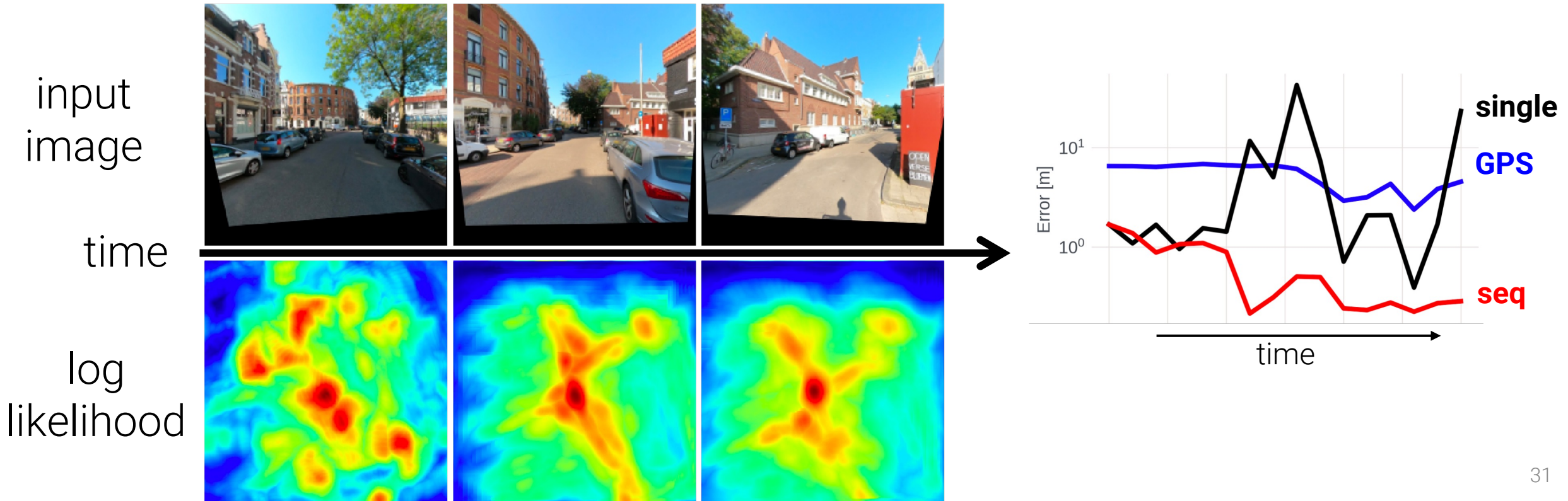
input
image

time →

log
likelihood

# Sequence localization

**Fuse successive predictions** assuming known relative poses
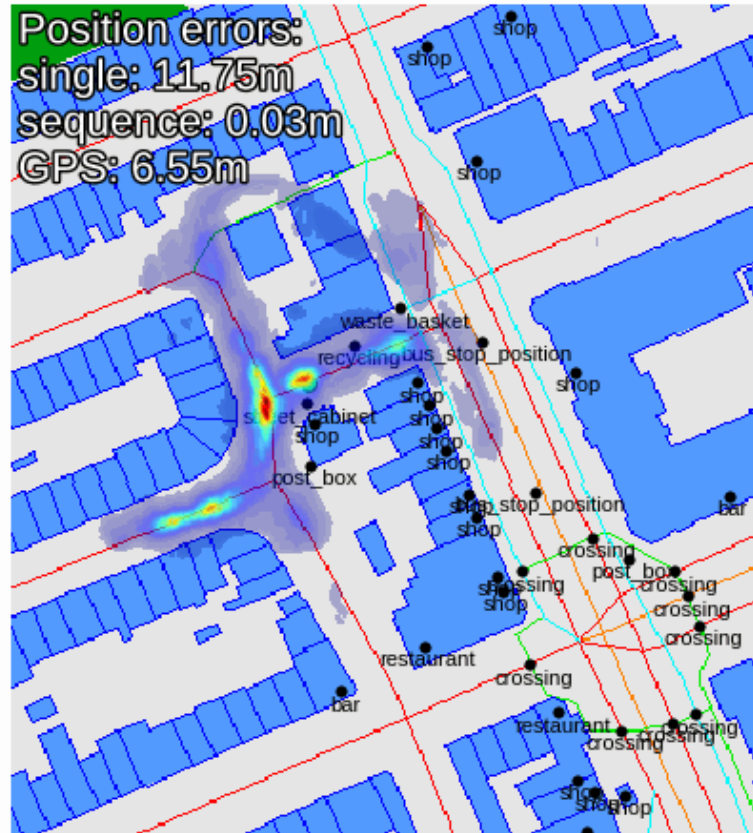
$$P(\boldsymbol{\xi}_i | \{\mathbf{I}_j\}, \mathrm{map}) = \prod_k P(\boldsymbol{\xi}_i \oplus \hat{\boldsymbol{\xi}}_{ij} | \mathbf{I}_j, \mathrm{map})$$

input image

time
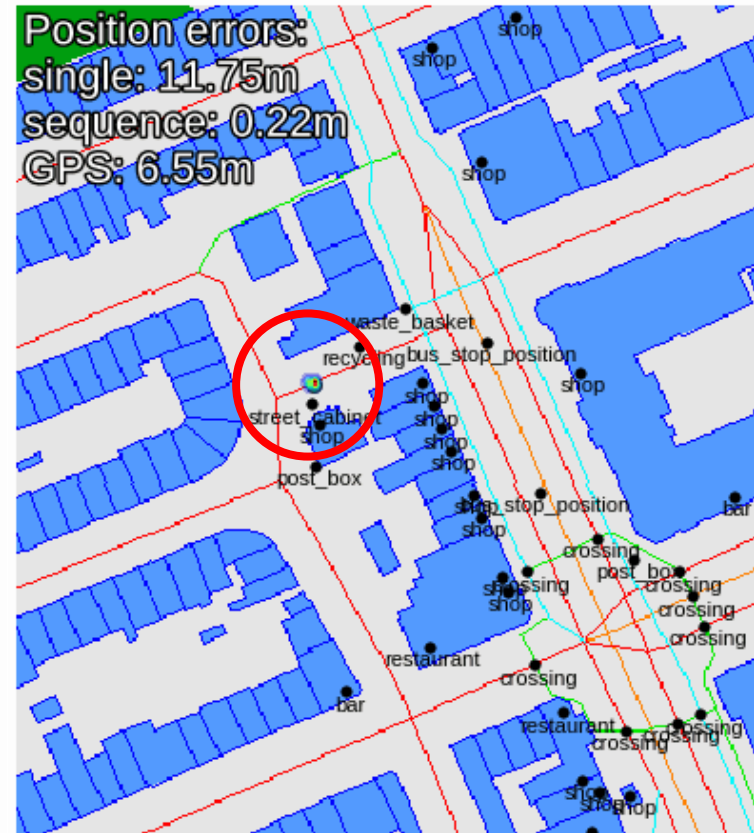
log likelihood

# Sequence localization



input image

single-frame likelihood

sequence likelihood

Position errors:
single: 11.75m
sequence: 0.03m
GPS: 6.55m

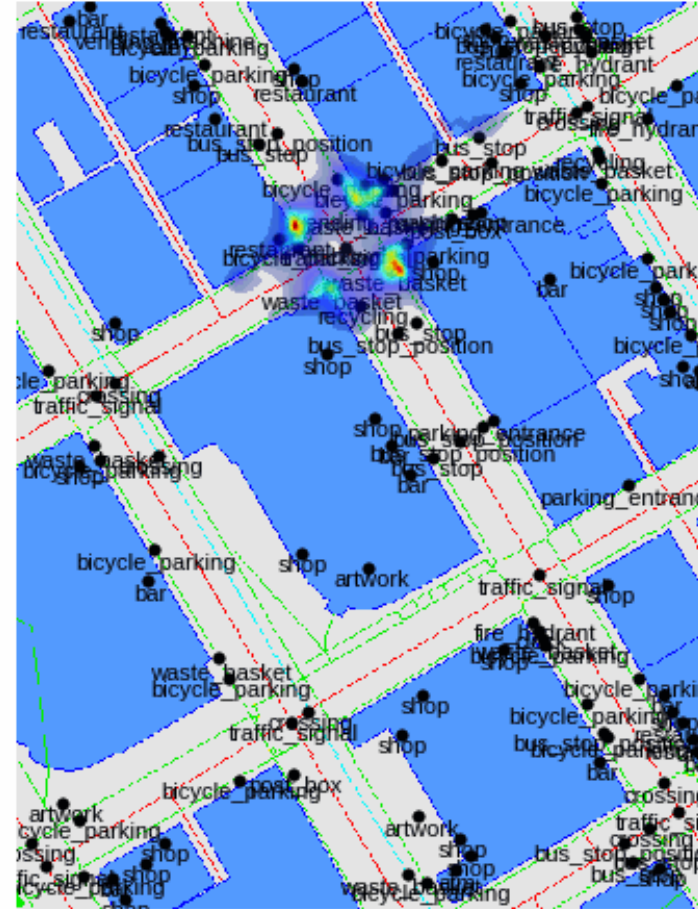Position errors:
single: 11.75m
sequence: 0.22m
GPS: 6.55m

# Sequence localization – Aria



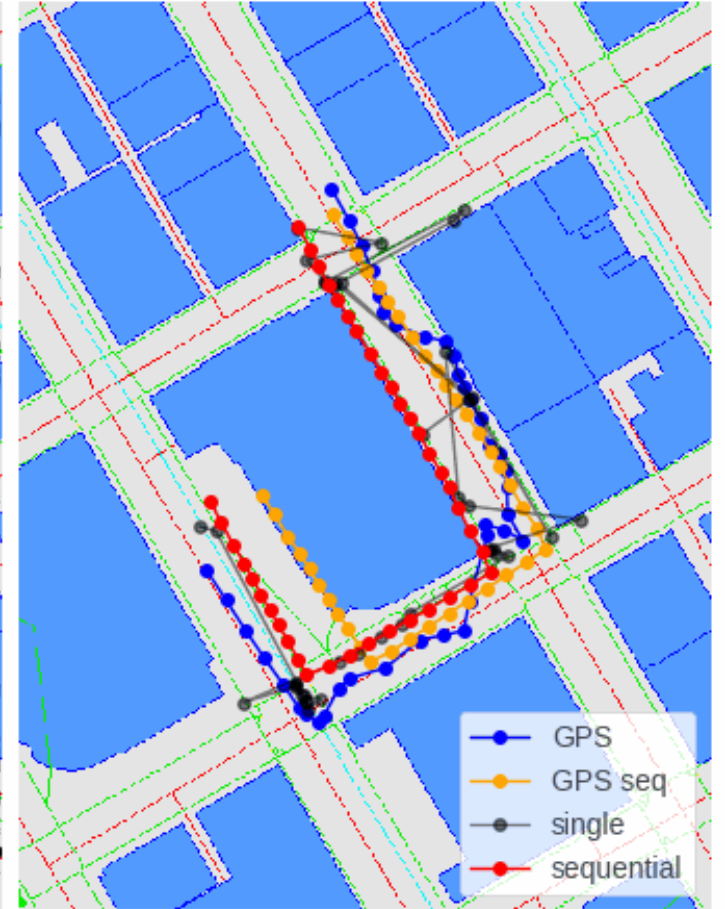input image      single-frame likelihood      final trajectories

# OrienterNet

## Visual Localization in 2D Public Maps with Neural Matching

Paul-Edouard Sarlin[1]    Daniel DeTone[2]    Tsun-Yi Yang[2]    Armen Avetisyan[2]

Julian Straub[2]    Tomasz Malisiewicz[2]    Samuel Rota Bulo[2]

Richard Newcombe[2]    Peter Kontschieder[2]    Vasileios Balntas[2]

CVPR 2023        psarlin.com/orienternet        Poster THU-PM-098