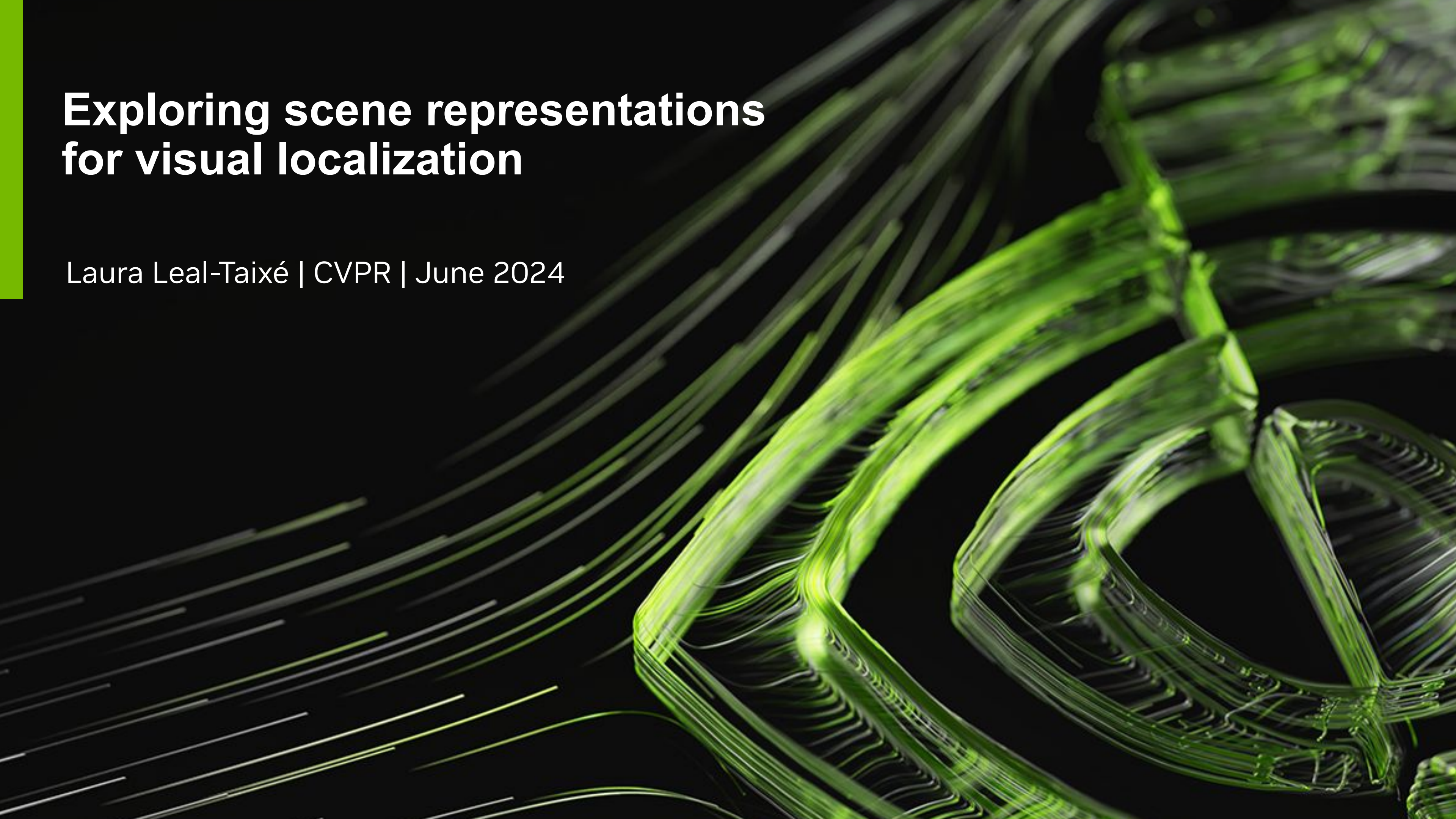# Exploring scene representations for visual localization
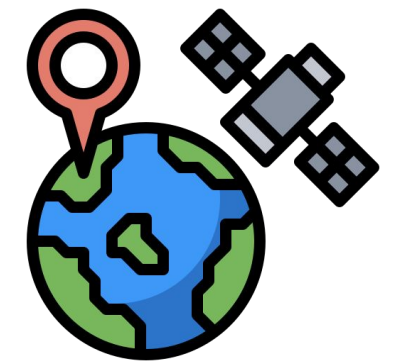
Laura Leal-Taixé | CVPR | June 2024
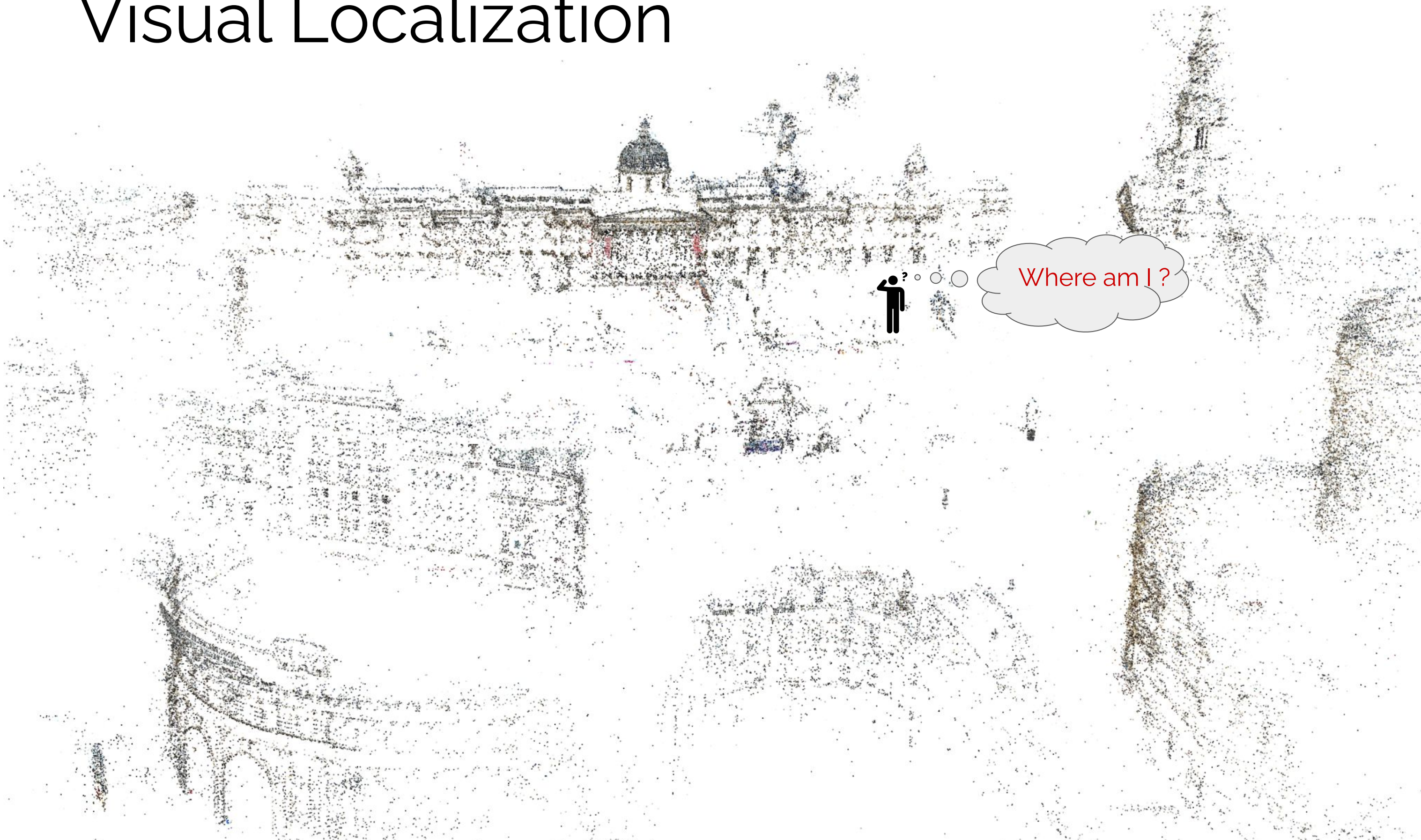
# Localization

Global Positioning
System
(GPS)

Where am I ?

# Visual Localization



Where am I ?
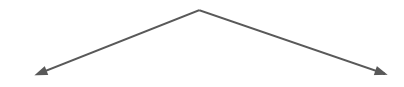
Localization System

Camera Pose

Orientation          Position (~cm)

(latitude, longitude) ~m

# Applications

## Indoor / Outdoor Navigation
(GPS-unavailable /unreliable)

## Autonomous Service Robots

## AR / VR
(Require cm-mm accuracy)

# Localization System

| Query | Scene Map | Method | Outputs |
|---|---|---|---|



Reference Images



Query Image



3D Point Cloud

# Structure-based Localization



Localize

2D/3D
Keypoints

Visual
Descriptors

2D-3D
Matching

Perspective-n-Point
Solver

# Practical Challenges

Storage Demand

| MegaDepth (192 scenes) | Camera | 3D Points | Images | Point Descriptors | | |
|---|---|---|---|---|---|---|
| | | | | SIFT | CAPS | SuperPoint |
| Storage | 15.73 MB | 3.44 GB | 157.84 GB | 130.10 GB | 520.38 GB | 1.041 TB |

# Practical Challenges

Storage Demand

| MegaDepth (192 scenes) | Camera | 3D Points | Images | Point Descriptors | | |
|---|---|---|---|---|---|---|
| | | | | SIFT | CAPS | SuperPoint |
| Storage | 15.73 MB | 3.44 GB | 157.84 GB | 130.10 GB | 520.38 GB | 1.041 TB |



3.44 GB

Scene Compression [1]

Desc Quantization [1, 2]

Scene Descriptor

130 GB ~ 1.04 TB

[1] Camposeco, Federico, et al. "Hybrid scene compression for visual localization." CVPR19
[2] Sattler, Torsten, Bastian Leibe, and Leif Kobbelt. "Efficient & effective prioritized matching for large-scale image-based localization." PAMI16

# Practical Challenges



Privacy Risk

Descriptor Inversion [1]

Client

Matching + Pose Estimation

Server

Matching + Pose Estimation

[1] Francesco, Pittaluga, et al Revealing Scenes by Inverting Structure From Motion Reconstructions. CVPR19

# Practical Challenges



Privacy Risk

Descriptor Inversion [1]

Client

Matching + Pose Estimation

Server

Matching + Pose Estimation

Privacy-preserving Descriptors [2, 3]

*Keypoints*  *Subspaces*  *Inversion*  *Reconstruction*

NinjaNet

input image  original descriptors  content-concealing descriptors

[1] Francesco, Pittaluga, et al Revealing Scenes by Inverting Structure From Motion Reconstructions. CVPR19
[2] Dusmanu, Mihai, et al. "Privacy-preserving image features via adversarial affine subspace embeddings." CVPR21.
[3] Ng, Tony, et al. "NinjaDesc: Content-Concealing Visual Descriptors via Adversarial Learning." CVPR22

# Practical Challenges



Maintenance
Complexity

Descriptor Upgrade [1]

SIFT        CAPS        SuperPoint

Cross-device Matching and Localization [1]

SIFT

SuperPoint

CAPS

SOSNet

HardNet

[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Practical Challenges



Maintenance Complexity

Descriptor Upgrade [1]

SIFT          CAPS          SuperPoint

Map Re-building

Upgraded Scene Descriptor

SIFT    software update    SOSNet        SOSNet    Translated SIFT    SIFT

matching                                    matching

SIFT    Translated SIFT                     SOSNet

Cross-device Matching and Localization [1]

SIFT

SuperPoint                    CAPS

SOSNet                    HardNet

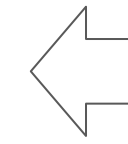[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Practical Challenges



**Storage Demand**

Scene Descriptor

Scene Compression

Desc Quantization

**Privacy Risk**

Client

Server

Privacy-preserving Descriptors

Keypoints    Subspaces    original descriptors    content-concealing descriptors

**Maintenance Complexity**

SIFT    CAPS    SuperPoint

SuperPoint    SIFT    CAPS    HardNet    SOSNet

Descriptors Translation

SIFT    software update    SOSNet    SOSNet    Translated SIFT    SIFT

matching    matching

SIFT    Translated SIFT    SOSNet

# Geometric-based Matching



Query

Keypoint Detection

Geometric-based Matching

Image Retrieval

2D/3D Keypoints

Visual Descriptor

Point Cloud Descriptor

2D-3D Matching

# Geometric-based Matching



Storage Demand

Minimal Representation

Privacy Risk

Client

Server

Maintenance Complexity

Scalable Large-scale Localization

🙂 Low Storage

🙂 Privacy Preserving

🙂 No Descriptor Maintenance

# Geometric-based matching and pose estimation

BPnPNet [4]
- Learning-based
- Declarative layers
- Degrades with outliers.

# Geometric-based matching and pose estimation

Does not scale to real-world localization settings!

BPnPNet [4]
- Learning-based
- Declarative layers
- Degrades with outliers.

# Geometric-based matching and pose estimation



GoMatch

Query · 2D Keypoint Pixels · Attention · Geometric Descriptors · Sinkhorn · Outlier Rejection

DB Keypoint Representations · 3D Coordinates · Bearing Vectors · shared if both BVs · E · E

BPnPNet [4]
- Learning-based
- Declarative layers
- Degrades with outliers.

2D Keypoints Query

3D Keypoints Retrieval

2D Feature Extractor

3D Feature Extractor

Sinkhorn Layer

Pose Solver

BPnPNet

Median Translation Error

Median Rotation Error (°)

# GoMatch Step-by-Step



2D Keypoints Query

2D Feature Extractor

Sinkhorn Layer

Pose Solver

2D Keypoints Retrieval

2D Feature Extractor

BPnPNet

BVs

Bearing Vectors instead of 3D keypoints

Median Translation Error

Median Rotation Error (°)

A bearing vector represents the direction from the reference camera origin to a 3D point in normalized coordinates.

2D Keypoints
Query

2D Feature
Extractor

2D Keypoints
Retrieval

2D Feature
Extractor

Sinkhorn
Layer

Pose Solver

Median Translation Error

BPnPNet
BVs

Shared Encoder

Shared feature
encoder

Median Rotation Error (°)

2D Keypoints Query

2D Feature Extractor

2D Keypoints Retrieval

2D Feature Extractor

Attention

Sinkhorn Layer

Outlier Rejection

Median Translation Error

Median Rotation Error (°)

BPnPNet

BVs

Shared Encoder

Attention
Outlier Rejection

Further removing outliers

2D Keypoints Query

2D Feature Extractor

Attention

Sinkhorn Layer

Outlier Rejection

2D Keypoints Retrieval

2D Feature Extractor

BPnPNet

BVs

Shared Encoder

Attention

Outlier Rejection

More covisible views

Increasing the number of retrieved images

Median Translation Error

Median Rotation Error (°)

# Generalization: outdoor/indoor and keypoints

MegaDepth
(Outdoor w.SIFT)



Rialto Bridge, Venice    Eiffel Tower, Paris    Central Park, NYC

Grand Canal, Venice    Trafalgar Square, London    Colosseum, Rome

# Generalization: outdoor/indoor and keypoints



MegaDepth
(Outdoor w.SIFT)

Rialto Bridge, Venice

Eiffel Tower, Paris

Central Park, NYC

Grand Canal, Venice

Trafalgar Square, London

Colosseum, Rome

Indoor

7Scenes

0.22

Median Translation Err (m)

GoMatch (w.SIFT)

7Scenes

5.78

Median Rotation Err (deg)

GoMatch (w.SIFT)

# Generalization: outdoor/indoor and keypoints

# Comparison with SOTA – Cambridge Landmarks



Median Translation Error (m)

Median Rotation Error (deg)

# Comparison with SOTA – Cambridge Landmarks



Median Translation Error (m)

Legend: GM, E2E, VM

HLoc 0.18, ActiveSearch 0.29, HybridSC 0.56, DSAC++ 0.14, MS-Transformer 1.28, MS-PoseNet 2.74, PoseNet 2.09, GoMatch 1.73, BPnPNet 11.44

Median Rotation Error (deg)

Legend: GM, E2E, VM

HLoc 0.63, ActiveSearch 0.63, HybridSC 0.66, DSAC++ 0.33, MS-Transformer 2.73, MS-PoseNet 5.34, PoseNet 6.84, GoMatch 11.02, BPnPNet 106.72

# Comparison with SOTA – Cambridge Landmarks

# Comparison with SOTA – 7 Scenes



Median Translation Err (m) — GM, E2E, VM. Values: HLoc 0.04, ActiveSearch 0.05, DSAC++ 0.04, MS-PoseNet 0.2, MS-Transformer 0.18, PoseNet 0.44, GoMatch 0.18, BPnPNet 1.61.

Median Rotation Err (deg) — GM, E2E, VM. Values: HLoc 1.59, ActiveSearch 2.46, DSAC++ 1.1, MS-PoseNet 8.41, MS-Transformer 7.28, PoseNet 10.4, GoMatch 4.61, BPnPNet 39.3.

Storage (MB) — GM, E2E, VM. Values: HLoc 22977, DSAC++ 1449, MS-Transformer 71.1, PoseNet 350, GoMatch 397, BPnPNet 397.

# Compare to VM – Cambridge Landmarks

| | Method | Storage (MB) | No Desc. Maint. | Privacy | King's College | Old Hospital Median Pose | Shop Facade Error (m, °) (↓) | St. Mary's Church |
|---|---|---|---|---|---|---|---|---|
| **E2E** | PoseNet [38] | 200 | ✓ | ✓ | 1.92/5.40 | 2.31/5.38 | 1.46/8.08 | 2.65/8.48 |
| | DSAC++ [6] | 828 | ✓ | ✓ | 0.18/0.30 | 0.20/0.30 | 0.06/0.30 | 0.13/0.40 |
| | MSPN [4] | - | ✓ | ✓ | 1.73/3.65 | 2.55/4.05 | 2.92/7.49 | 2.67/6.18 |
| | MS-Transformer [65] | 71.1 | ✓ | ✓ | 0.83/1.47 | 1.81/2.39 | 0.86/3.07 | 1.62/3.99 |
| **VM** | HybridSC [14] | 3.13 | ✗ | ? | 0.81/0.59 | 0.75/1.01 | 0.19/0.54 | 0.50/0.49 |
| | Active Search [58] | 812.7 | ✗ | ✗ | 0.42/0.55 | 0.44/1.01 | 0.12/0.40 | 0.19/0.54 |
| | HLoc [55](w.SP [22]) | 3214.84 | ✗ | ✗ | 0.16/0.38 | 0.33/1.04 | 0.07/0.54 | 0.16/0.54 |
| | HLoc(w.SP+SG [56]) | 3214.84 | ✗ | ✗ | **0.12/0.20** | **0.15/0.30** | **0.04/0.20** | **0.07/0.21** |
| **GM** | BPnPNet [11] | 48.15 | ✓ | ✓ | 26.73/106.99 | 24.8/162.99 | 7.53/107.17 | 11.11/49.74 |
| | GoMatch | 48.15 | ✓ | ✓ | 0.25/0.64 | 2.83/8.14 | 0.48/4.77 | 3.35/9.94 |

# Compare to VM – Cambridge Landmarks



| | Method | Storage (MB) | No Desc. Maint. | Privacy | King's College | Old Hospital Median Pose | Shop Facade Error (m, °) (↓) | St. Mary's Church |
|---|---|---|---|---|---|---|---|---|
| E2E | PoseNet [38] | 200 | ✓ | ✓ | 1.92/5.40 | 2.31/5.38 | 1.46/8.08 | 2.65/8.48 |
| | DSAC++ [6] | 828 | ✓ | ✓ | 0.18/0.30 | 0.20/0.30 | 0.06/0.30 | 0.13/0.40 |
| | MSPN [4] | - | ✓ | ✓ | 1.73/3.65 | 2.55/4.05 | 2.92/7.49 | 2.67/6.18 |
| | MS-Transformer [65] | 71.1 | ✓ | ✓ | 0.83/1.47 | 1.81/2.39 | 0.86/3.07 | 1.62/3.99 |
| VM | HybridSC [14] | 3.13 | ✗ | ? | 0.81/0.59 | 0.75/1.01 | 0.19/0.54 | 0.50/0.49 |
| | Active Search [58] | 812.7 | ✗ | ✗ | 0.42/0.55 | 0.44/1.01 | 0.12/0.40 | 0.19/0.54 |
| | HLoc [55](w.SP [22]) | 3214.84 | ✗ | ✗ | 0.16/0.38 | 0.33/1.04 | 0.07/0.54 | 0.16/0.54 |
| | HLoc(w.SP+SG [56]) | 3214.84 | ✗ | ✗ | **0.12/0.20** | **0.15/0.30** | **0.04/0.20** | **0.07/0.21** |
| GM | BPnPNet [11] | 48.15 | ✓ | ✓ | 26.73/106.99 | 24.8/162.99 | 7.53/107.17 | 11.11/49.74 |
| | GoMatch | 48.15 | ✓ | ✓ | 0.25/0.64 | 2.83/8.14 | 0.48/4.77 | 3.35/9.94 |

# Compare to VM – Cambridge Landmarks



| | Method | Storage (MB) | No Desc. Maint. | Privacy | King's College | Old Hospital Median Pose | Shop Facade Error (m, °) (↓) | St. Mary's Church |
|---|---|---|---|---|---|---|---|---|
| **E2E** | PoseNet [38] | 200 | ✓ | ✓ | 1.92/5.40 | 2.31/5.38 | 1.46/8.08 | 2.65/8.48 |
| | DSAC++ [6] | 828 | ✓ | ✓ | 0.18/0.30 | 0.20/0.30 | 0.06/0.30 | 0.13/0.40 |
| | MSPN [4] | - | ✓ | ✓ | 1.73/3.65 | 2.55/4.05 | 2.92/7.49 | 2.67/6.18 |
| | MS-Transformer [65] | 71.1 | ✓ | ✓ | 0.83/1.47 | 1.81/2.39 | 0.86/3.07 | 1.62/3.99 |
| **VM** | HybridSC [14] | 3.13 | ✗ | ? | 0.81/0.59 | 0.75/1.01 | 0.19/0.54 | 0.50/0.49 |
| | Active Search [58] | 812.7 | ✗ | ✗ | 0.42/0.55 | 0.44/1.01 | 0.12/0.40 | 0.19/0.54 |
| | HLoc [55](w.SP [22]) | 3214.84 | ✗ | ✗ | 0.16/0.38 | 0.33/1.04 | 0.07/0.54 | 0.16/0.54 |
| | HLoc(w.SP+SG [56]) | 3214.84 | ✗ | ✗ | **0.12/0.20** | **0.15/0.30** | **0.04/0.20** | **0.07/0.21** |
| **GM** | BPnPNet [11] | 48.15 | ✓ | ✓ | 26.73/106.99 | 24.8/162.99 | 7.53/107.17 | 11.11/49.74 |
| | GoMatch | 48.15 | ✓ | ✓ | 0.25/0.64 | 2.83/8.14 | 0.48/4.77 | 3.35/9.94 |

# Compare to VM – Cambridge Landmarks

| | Method | Storage (MB) | No Desc. Maint. | Privacy | King's College | Old Hospital Median Pose | Shop Facade Error (m, °) (↓) | St. Mary's Church |
|---|---|---|---|---|---|---|---|---|
| E2E | PoseNet [38] | 200 | ✓ | ✓ | 1.92/5.40 | 2.31/5.38 | 1.46/8.08 | 2.65/8.48 |
| | DSAC++ [6] | 828 | ✓ | ✓ | 0.18/0.30 | 0.20/0.30 | 0.06/0.30 | 0.13/0.40 |
| | MSPN [4] | - | ✓ | ✓ | 1.73/3.65 | 2.55/4.05 | 2.92/7.49 | 2.67/6.18 |
| | MS-Transformer [65] | 71.1 | ✓ | ✓ | 0.83/1.47 | 1.81/2.39 | 0.86/3.07 | 1.62/3.99 |
| VM | HybridSC [14] | 3.13 | ✗ | ? | 0.81/0.59 | 0.75/1.01 | 0.19/0.54 | 0.50/0.49 |
| | Active Search [58] | 812.7 | ✗ | ✗ | 0.42/0.55 | 0.44/1.01 | 0.12/0.40 | 0.19/0.54 |
| | HLoc [55](w.SP [22]) | 3214.84 | ✗ | ✗ | 0.16/0.38 | 0.33/1.04 | 0.07/0.54 | 0.16/0.54 |
| | HLoc(w.SP+SG [56]) | 3214.84 | ✗ | ✗ | **0.12/0.20** | **0.15/0.30** | **0.04/0.20** | **0.07/0.21** |
| GM | BPnPNet [11] | 48.15 | ✓ | ✓ | 26.73/106.99 | 24.8/162.99 | 7.53/107.17 | 11.11/49.74 |
| | GoMatch | 48.15 | ✓ | ✓ | 0.25/0.64 | 2.83/8.14 | 0.48/4.77 | 3.35/9.94 |

# Compare to VM – Cambridge Landmarks

| | Method | Storage (MB) | No Desc. Maint. | Privacy | King's College | Old Hospital Median Pose | Shop Facade Error (m, °) (↓) | St. Mary's Church |
|---|---|---|---|---|---|---|---|---|
| E2E | PoseNet [38] | 200 | ✓ | ✓ | 1.92/5.40 | 2.31/5.38 | 1.46/8.08 | 2.65/8.48 |
| | DSAC++ [6] | 828 | ✓ | ✓ | 0.18/0.30 | 0.20/0.30 | 0.06/0.30 | 0.13/0.40 |
| | MSPN [4] | - | ✓ | ✓ | 1.73/3.65 | 2.55/4.05 | 2.92/7.49 | 2.67/6.18 |
| | MS-Transformer [65] | 71.1 | ✓ | ✓ | 0.83/1.47 | 1.81/2.39 | 0.86/3.07 | 1.62/3.99 |
| VM | HybridSC [14] | 3.13 | ✗ | ? | 0.81/0.59 | 0.75/1.01 | 0.19/0.54 | 0.50/0.49 |
| | Active Search [58] | 812.7 | ✗ | ✗ | 0.42/0.55 | 0.44/1.01 | 0.12/0.40 | 0.19/0.54 |
| | HLoc [55](w.SP [22]) | 3214.84 | ✗ | ✗ | 0.16/0.38 | 0.33/1.04 | 0.07/0.54 | 0.16/0.54 |
| | HLoc(w.SP+SG [56]) | 3214.84 | ✗ | ✗ | **0.12/0.20** | **0.15/0.30** | **0.04/0.20** | **0.07/0.21** |
| GM | BPnPNet [11] | 48.15 | ✓ | ✓ | 26.73/106.99 | 24.8/162.99 | 7.53/107.17 | 11.11/49.74 |
| | GoMatch | 48.15 | ✓ | ✓ | 0.25/0.64 | 2.83/8.14 | 0.48/4.77 | 3.35/9.94 |

# Conclusions

- Geometric localization is possible and (somewhat) SOTA

- Opens a new door for new work in privacy-aware, scalable localization

# Localization System

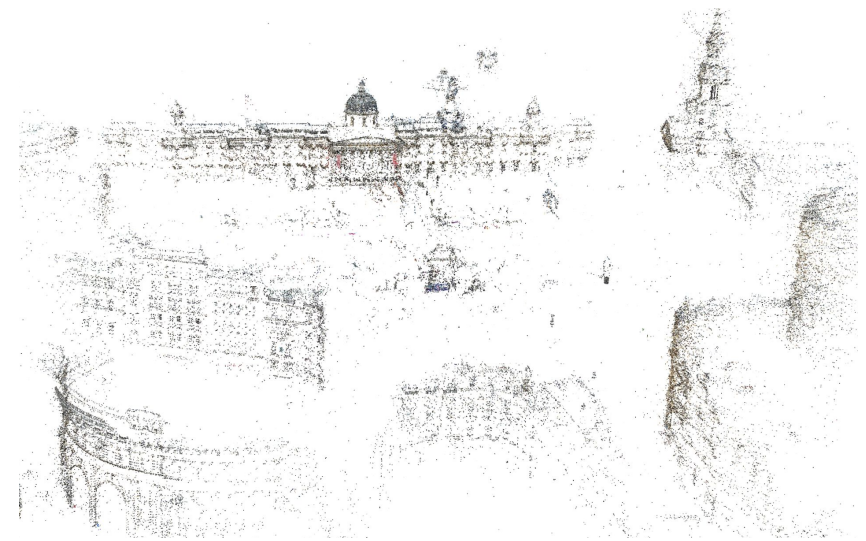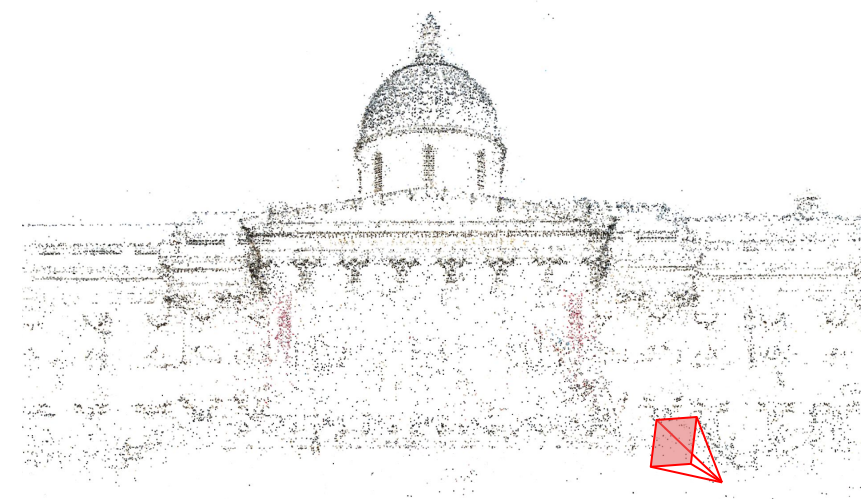| Query | + | Scene Map | ⇨ | Method | ⇨ | Outputs |
|---|---|---|---|---|---|---|



Reference Images



Query Image



**Descriptor-free**
Point Cloud

GoMatch



Even more compact
scene representation ?
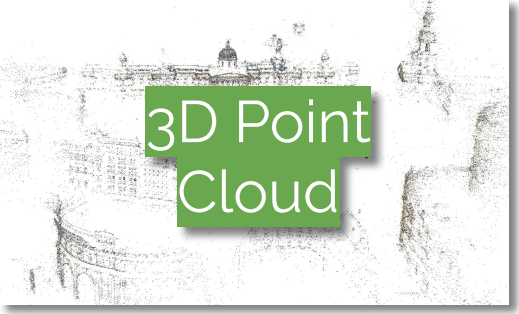
# Localization System



Query + Scene Map ⇨ Method ⇨ Outputs
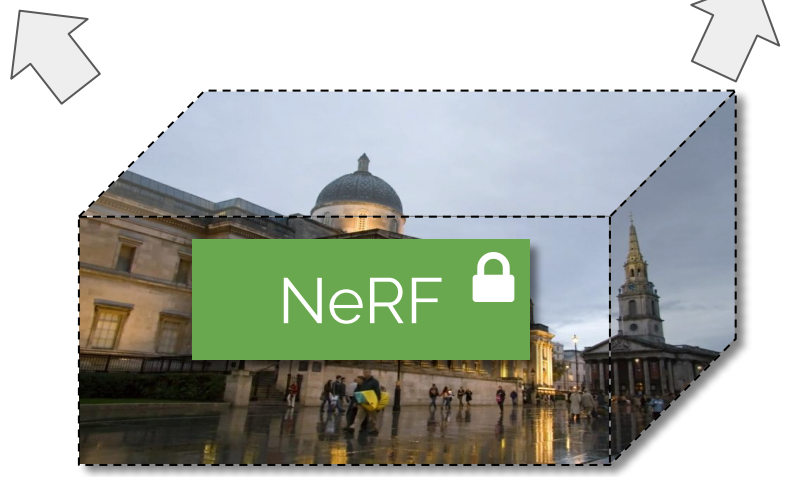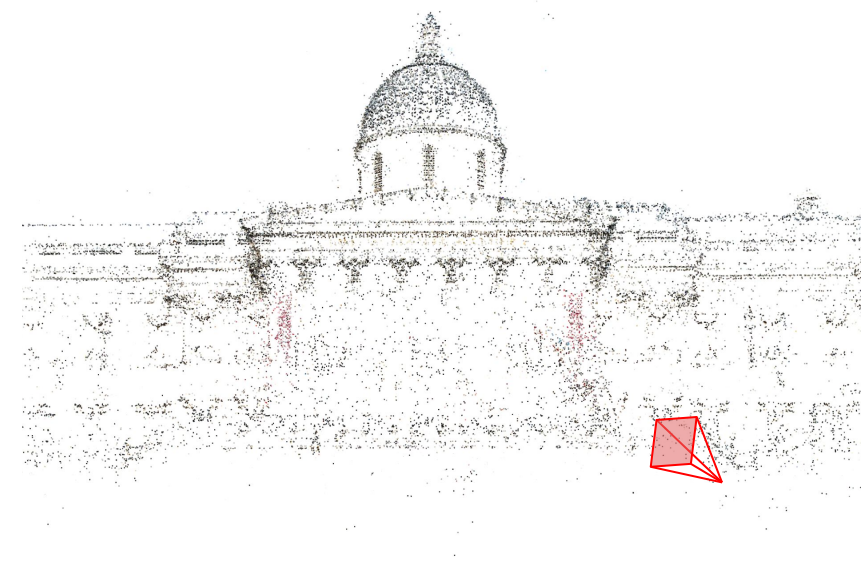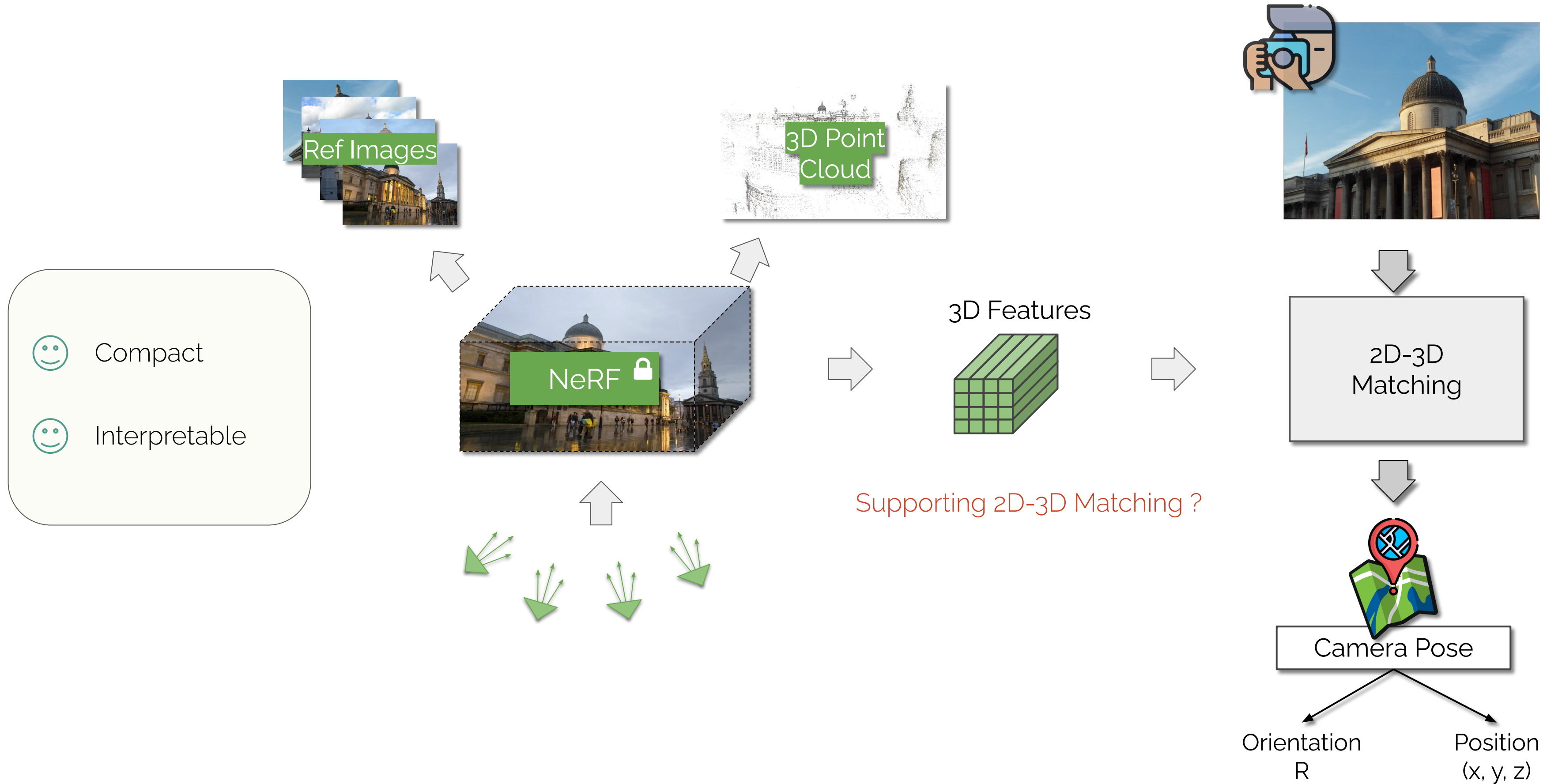
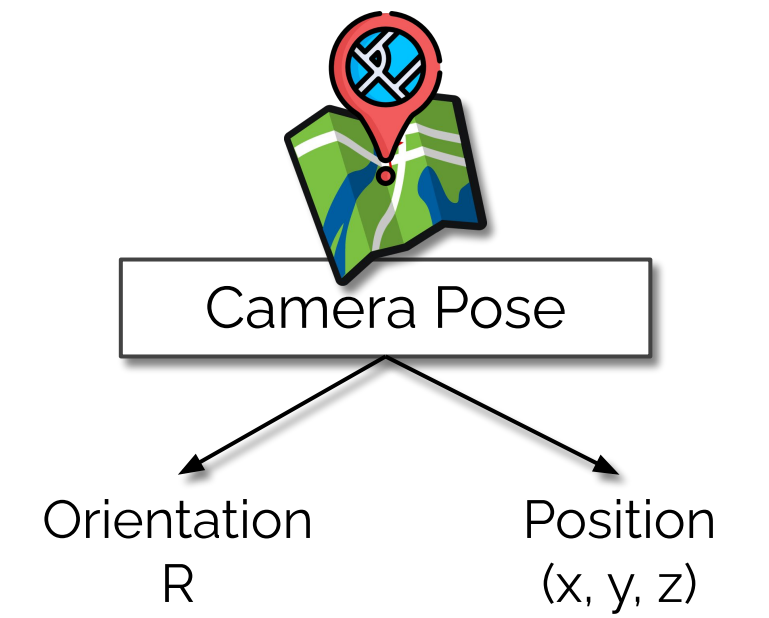Query Image

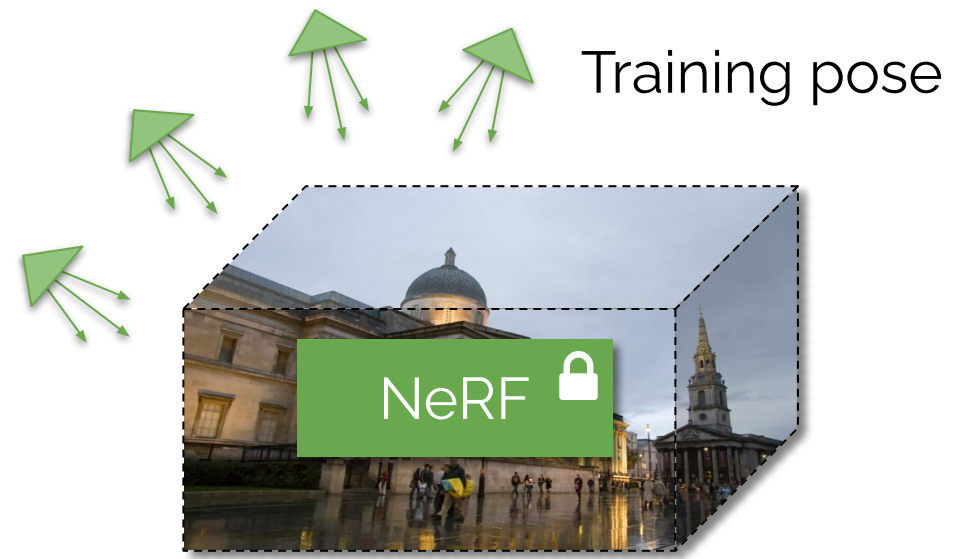Ref Images

3D Point Cloud

NeRF 🔒

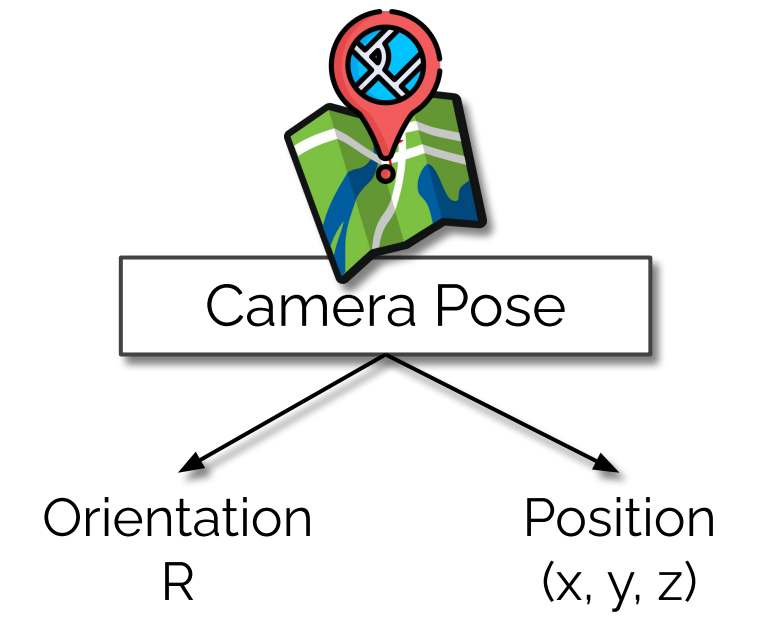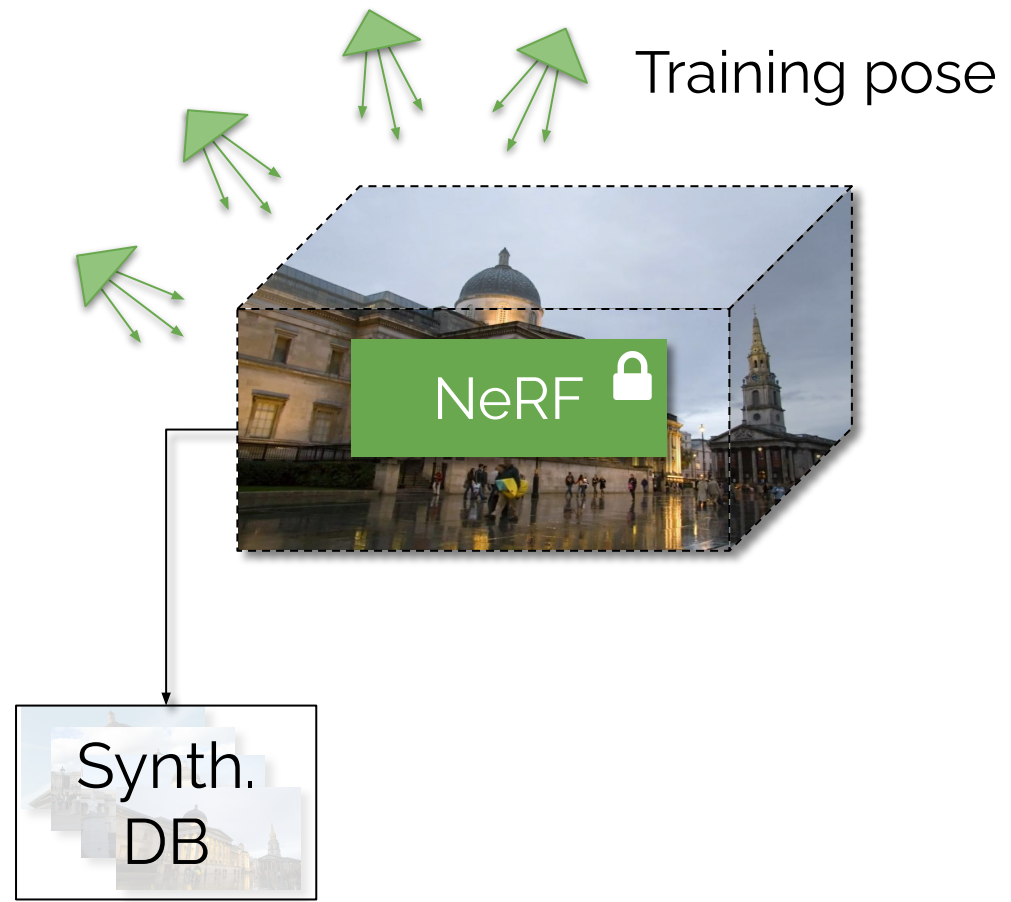NeRF-based Localization

☺ Compact

☺ Interpretable

# Introduction



Compact

Interpretable

Ref Images

3D Point Cloud

NeRF 🔒

3D Features

Supporting 2D-3D Matching ?

2D-3D Matching

Camera Pose

Orientation R

Position (x, y, z)

# Method



Training pose

NeRF 🔒

Camera Pose

Orientation
R

Position
(x, y, z)

Query Image

# Method



Training pose

NeRF 🔒

Synth.
DB

Camera Pose

Orientation
R

Position
(x, y, z)

Query Image

# Method



Training pose

NeRF 🔒

Synth. DB

Reference pose

Image Retrieval 🔒

Query Image

Camera Pose

Orientation R

Position (x, y, z)

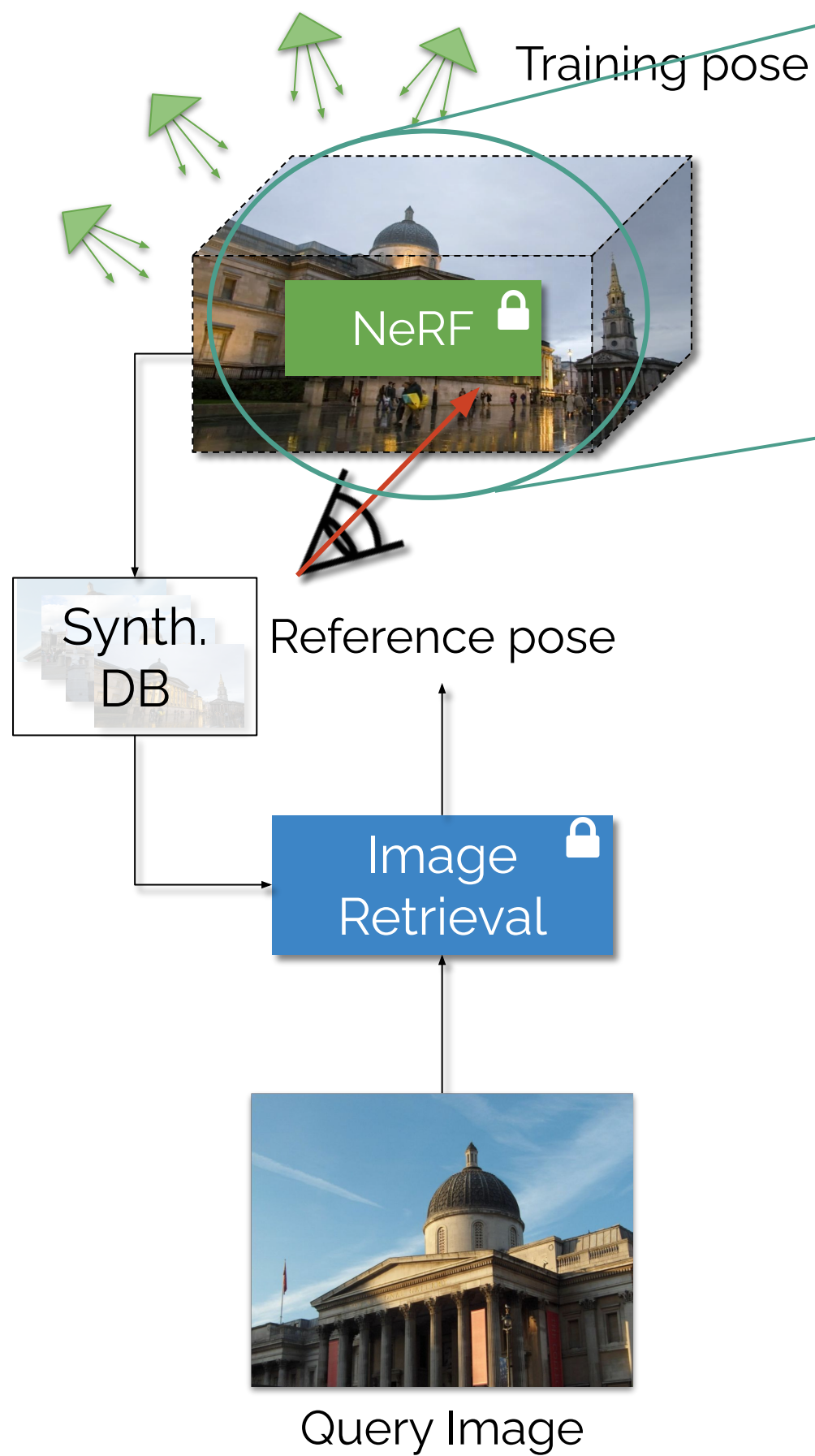# Method



Training pose

$$P_x(X) \longrightarrow \boxed{\Theta_x^1} \rightarrow \bullet\bullet\bullet \rightarrow \boxed{\Theta_x^j} \rightarrow \bullet\bullet\bullet \rightarrow \boxed{\Theta_x^L} \longrightarrow \boxed{\Theta_\sigma} \rightarrow \boldsymbol{\sigma}$$

NeRF

Rendering $f^j$
(Eq. 2)

$P_d(d) \longrightarrow \oplus \rightarrow \boxed{\Theta_c} \rightarrow \boldsymbol{c}$

NeRF Features $\hat{F}^j(r)$:

NeRF 🔒

Reference pose

Synth. DB

Image Retrieval 🔒

Query Image

Camera Pose

Orientation
R

Position
(x, y, z)

# Method



Training pose

3D (surface) Points

NeRF

NeRF Features

Synth. DB

Reference pose

Image Retrieval

Query Image

2D Encoder

Image Features

Camera Pose

Orientation
R

Position
(x, y, z)

# Method

# Method



3D (surface) Points

NeRF

Synth. DB

Reference pose

NeRF Features

Image Retrieval

NeRFMatch

2D-3D Matches

PnP Solver

Camera Pose

Orientation R

Position (x, y, z)

2D Encoder

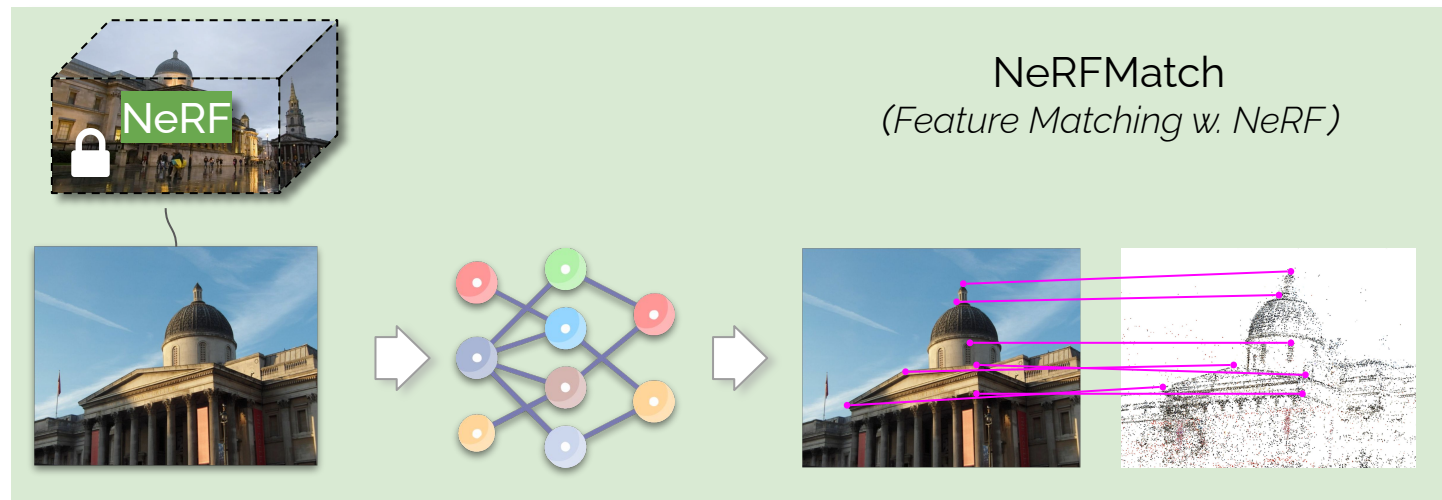Query Image
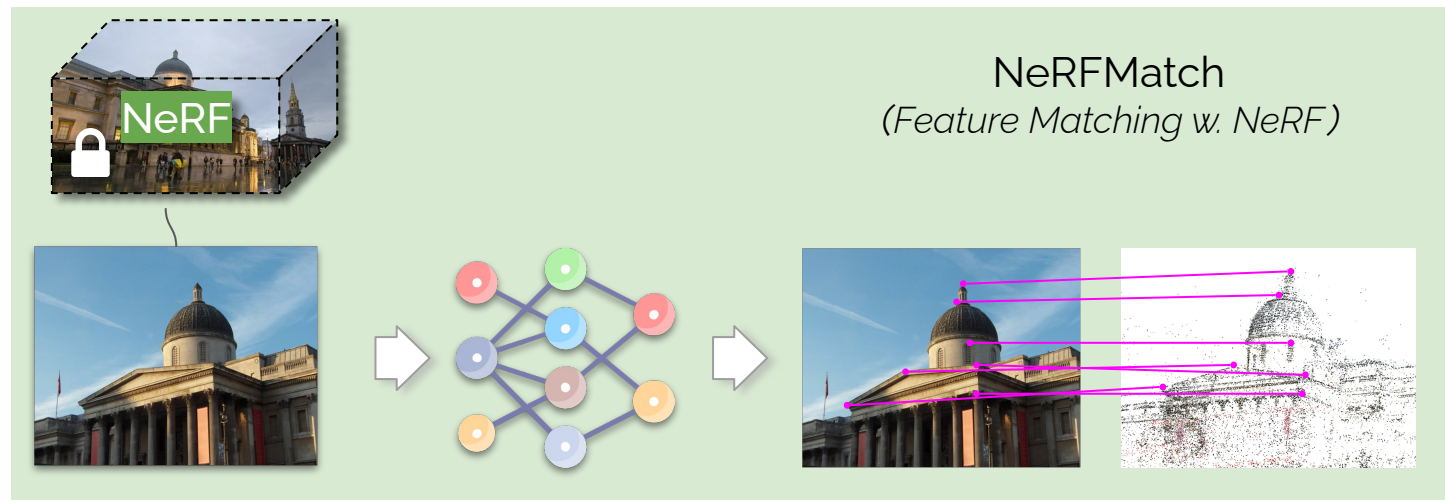
Image Features

NeRFMatch
*(Feature Matching w. NeRF)*

- NeRF not only provides 3D geometry but also comes with feature representation of 3D points that supports 2d-3D matching

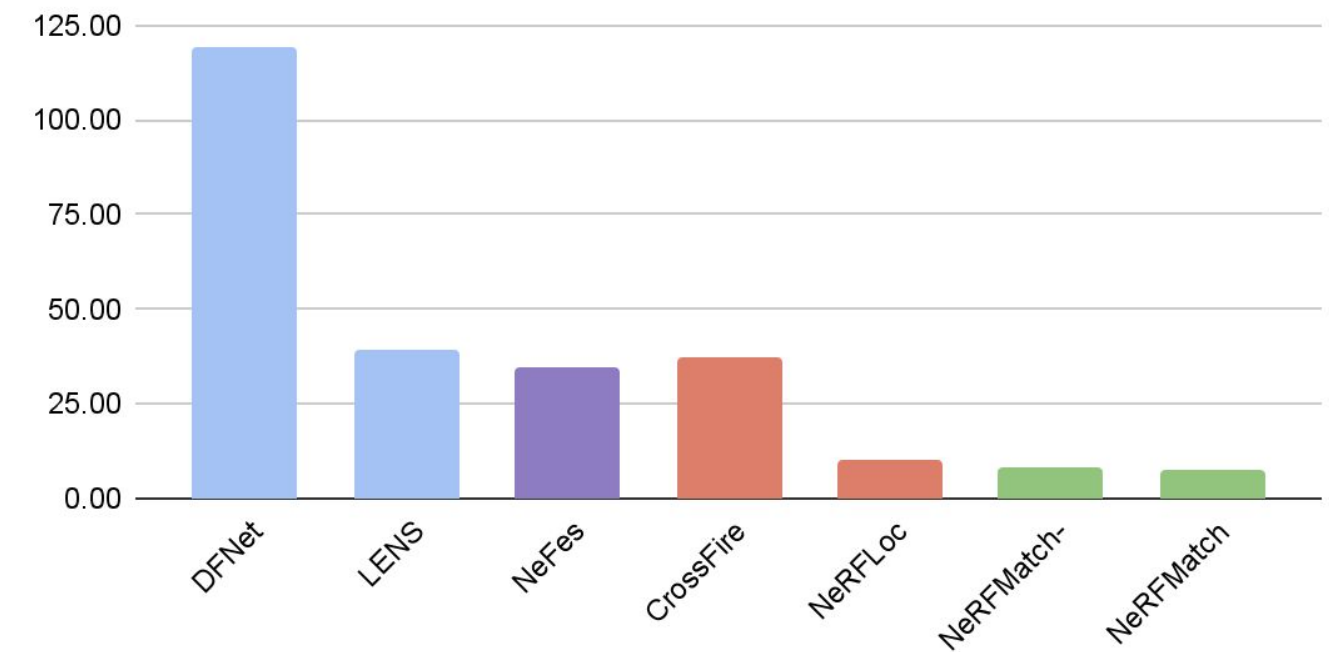| Metrics | Pt3D | Pe3D | $f^1$ | $f^2$ | $f^3$ | $f^4$ | $f^5$ | $f^6$ | $f^7$ |
|---|---|---|---|---|---|---|---|---|---|
| Med. Translation $(cm, \downarrow)$ | 432.3 | 25.5 | 22.9 | 23.3 | **21.8** | 22.3 | 23.5 | 24.1 | 40.9 |
| Med. Rotation $(°, \downarrow)$ | 6.5 | 0.6 | **0.5** | **0.5** | **0.5** | **0.5** | **0.5** | **0.5** | 1.0 |
| Localize Recall. $(\%, \uparrow)$ | 2.1 | 60.1 | 64.8 | 64.2 | **65.8** | 64.1 | 63.3 | 62.1 | 43.2 |

NeRFMatch
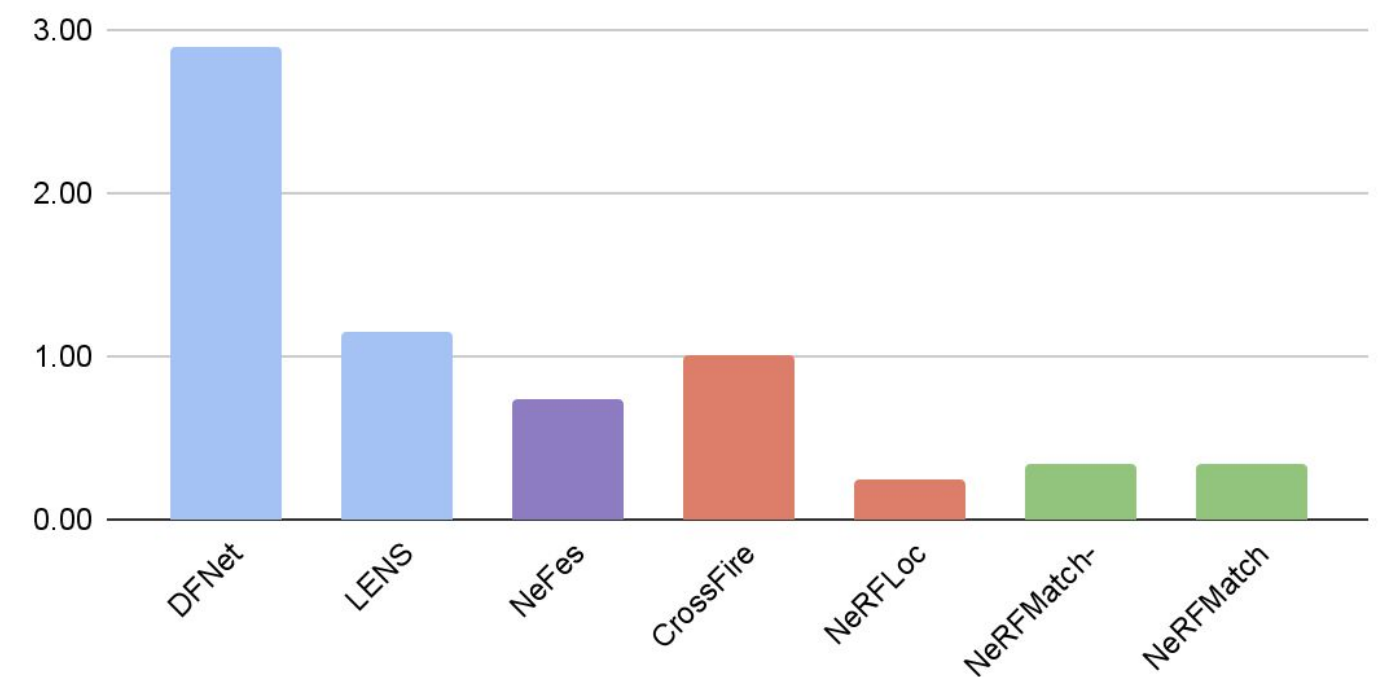*(Feature Matching w. NeRF)*
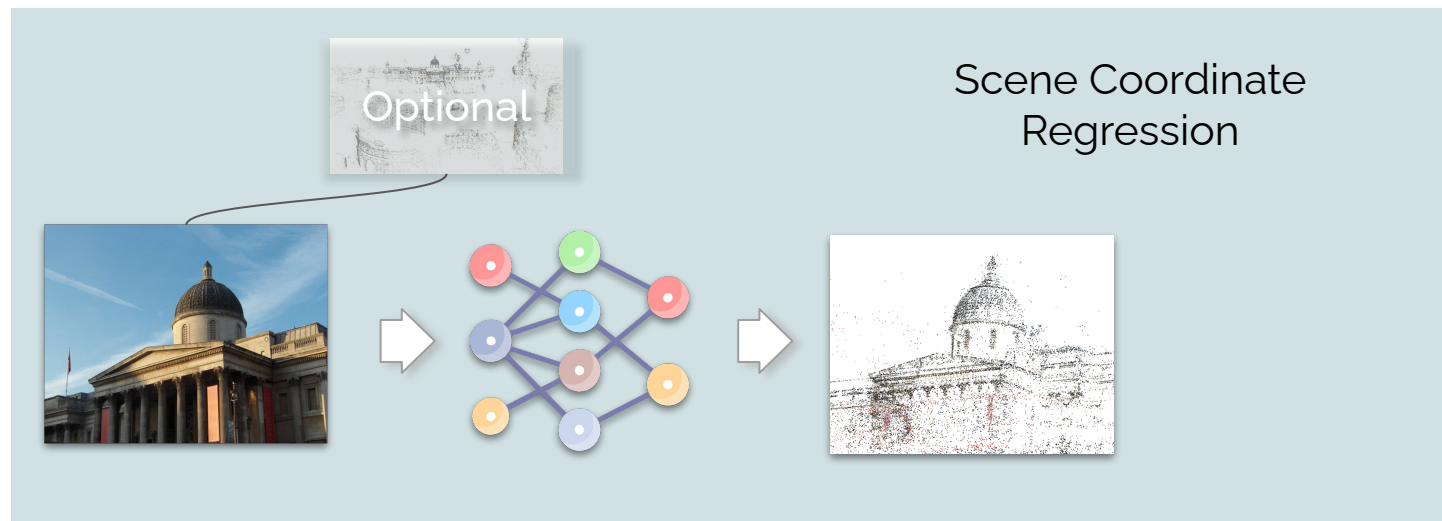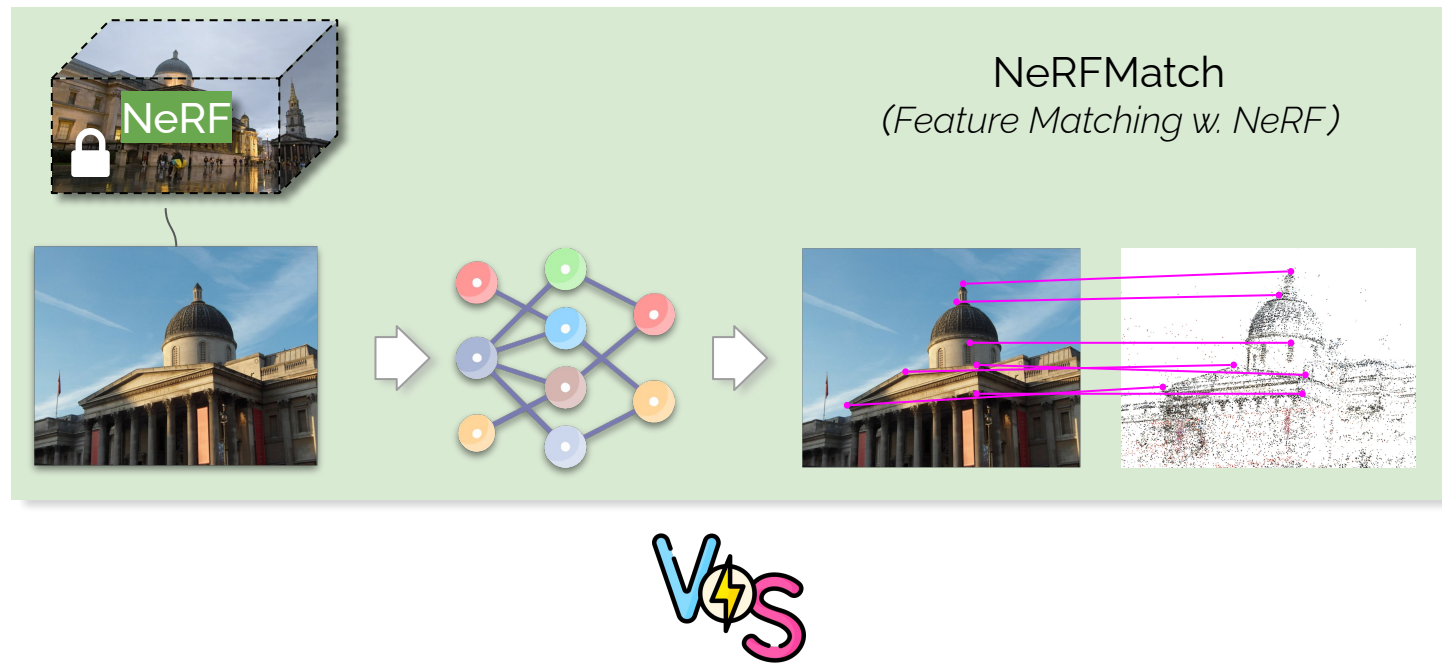
Compared to other ways of using NeRF:
- Training Augmentation: DFNet, LENS
- Test-time refinement: NeFes
- Joint NeRF and matching training: CrossFire, NeRFLoc (requires depth image input)
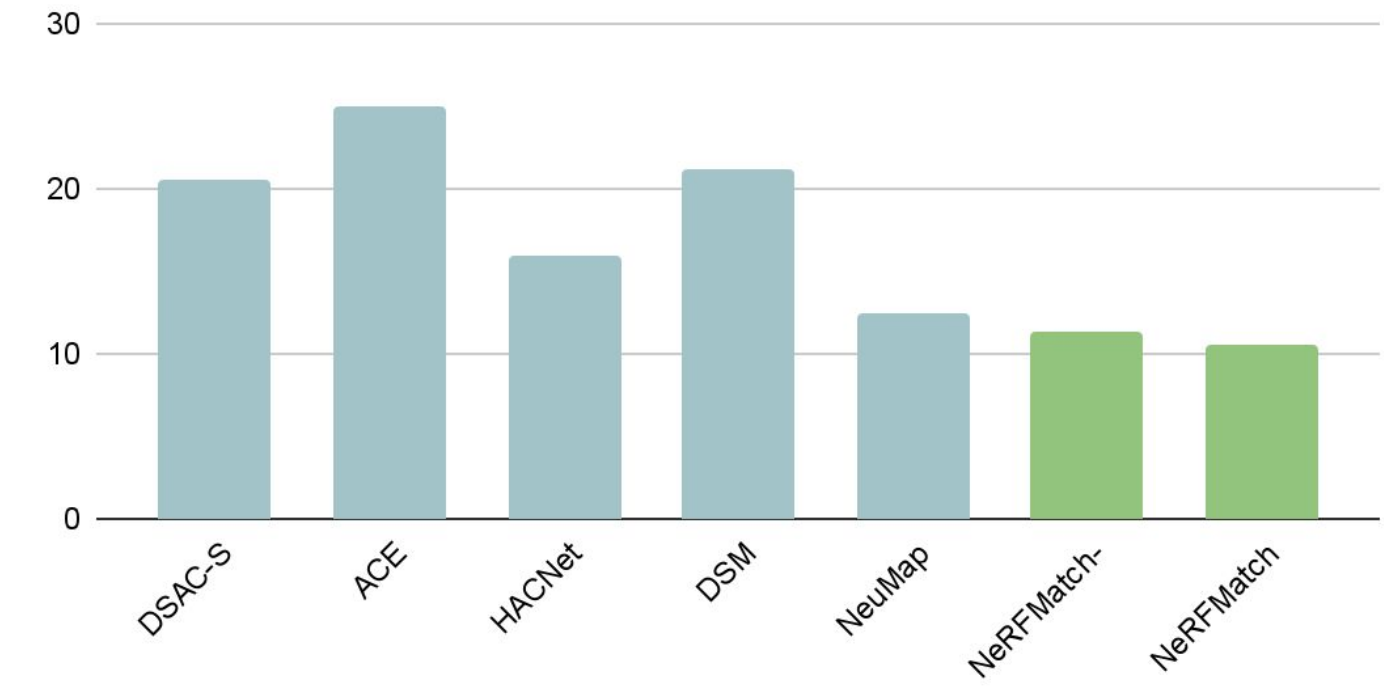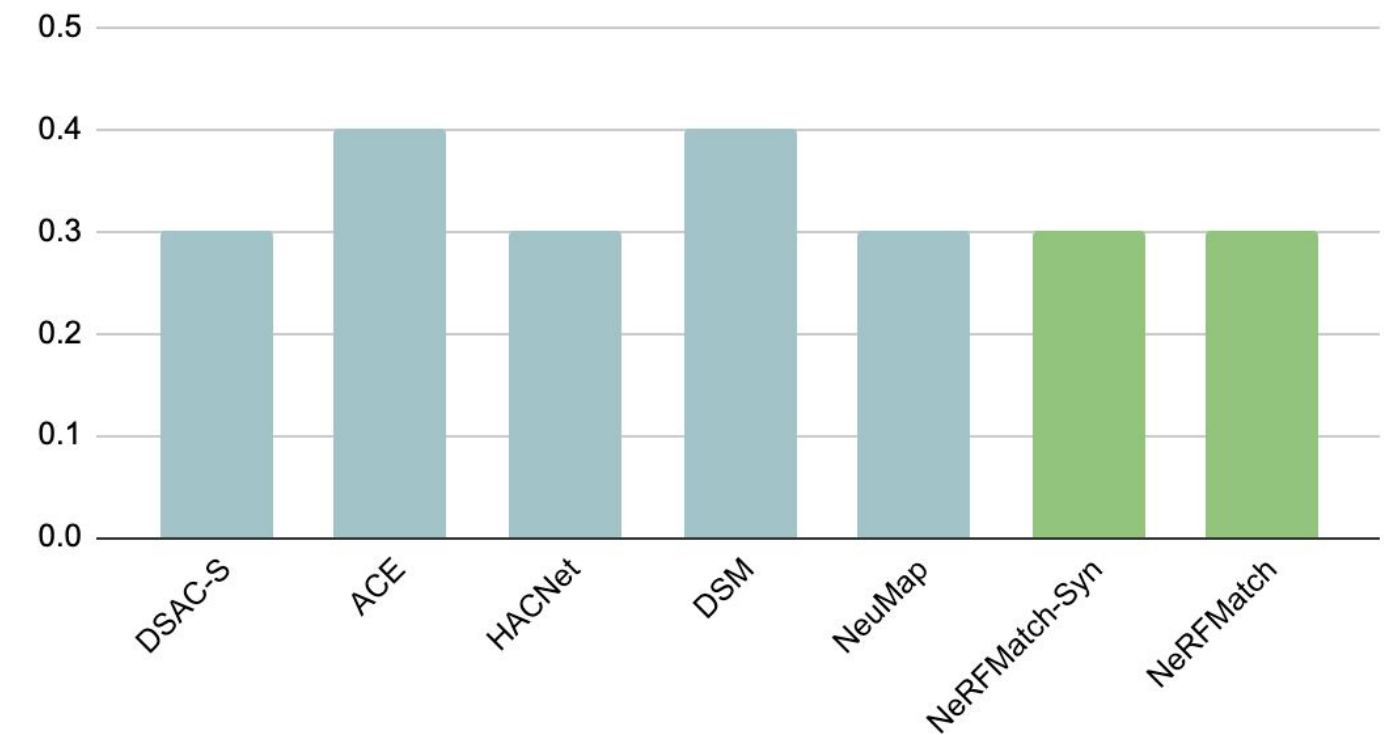


Cambridge Landmarks - Translation Error (cm)
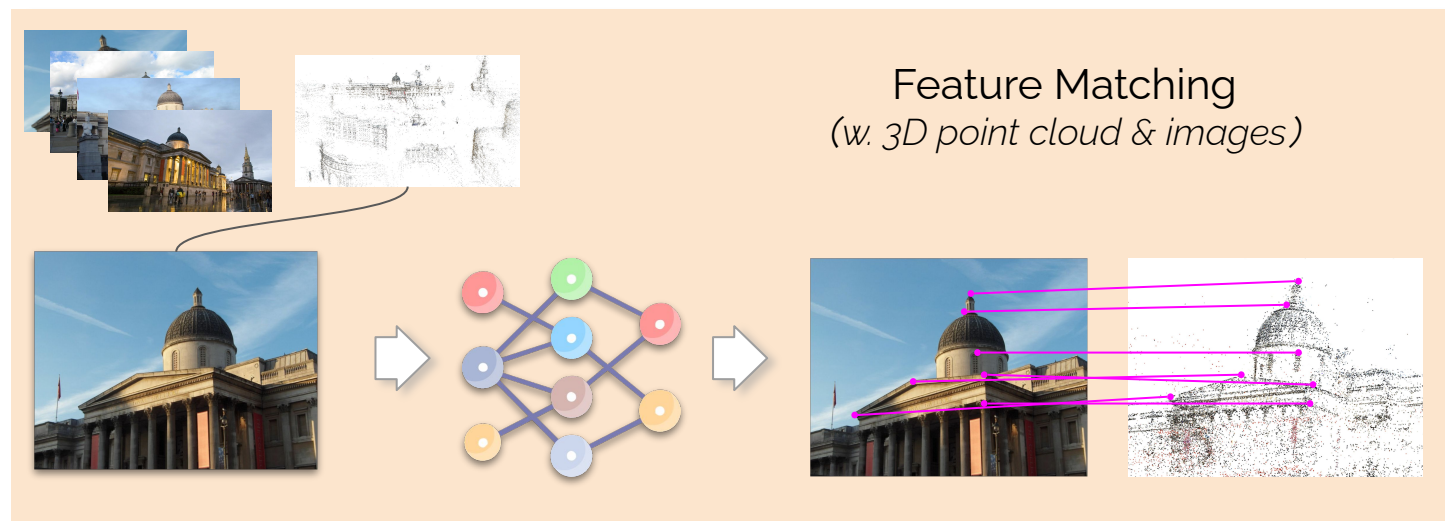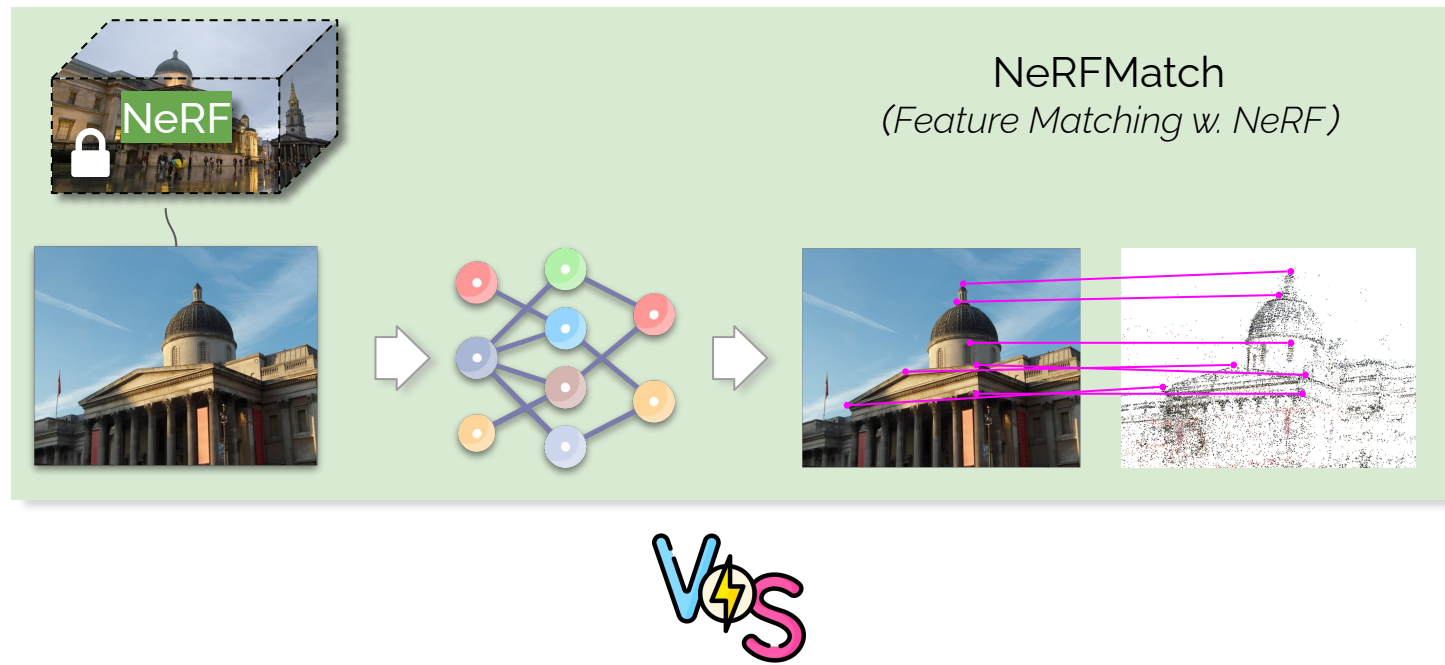


Cambridge Landmarks - Rotation Error (°)

NeRFMatch
*(Feature Matching w. NeRF)*
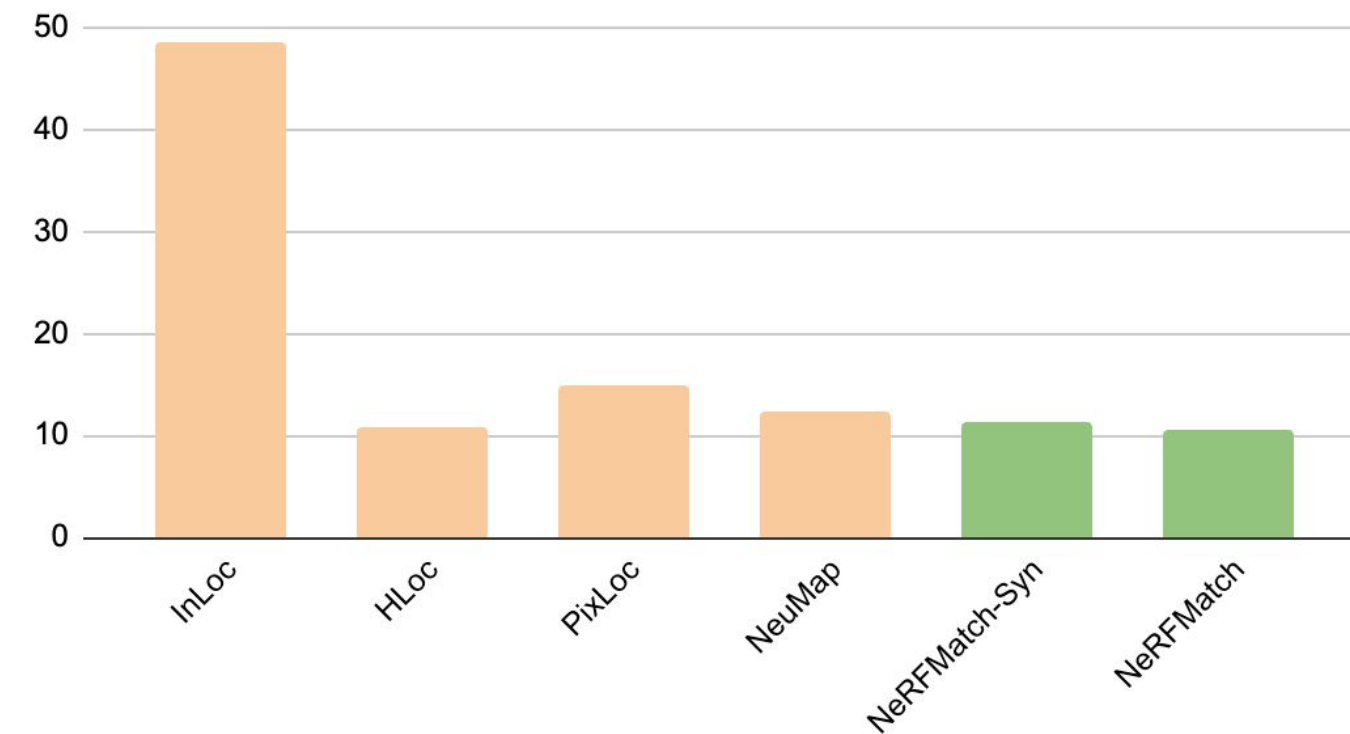
Scene Coordinate Regression

Optional

- Our method is similar to SCR methods where they map RGB to 3D points. Yet SCR represent 3D points with its coordiantes, while we use high-dimensional representation learned from NVS. And instead of regression, we learning a matching function to find a common ground between the (2d AND 3d) feature spaces. Our method outperform SCR on outdoor scenes.
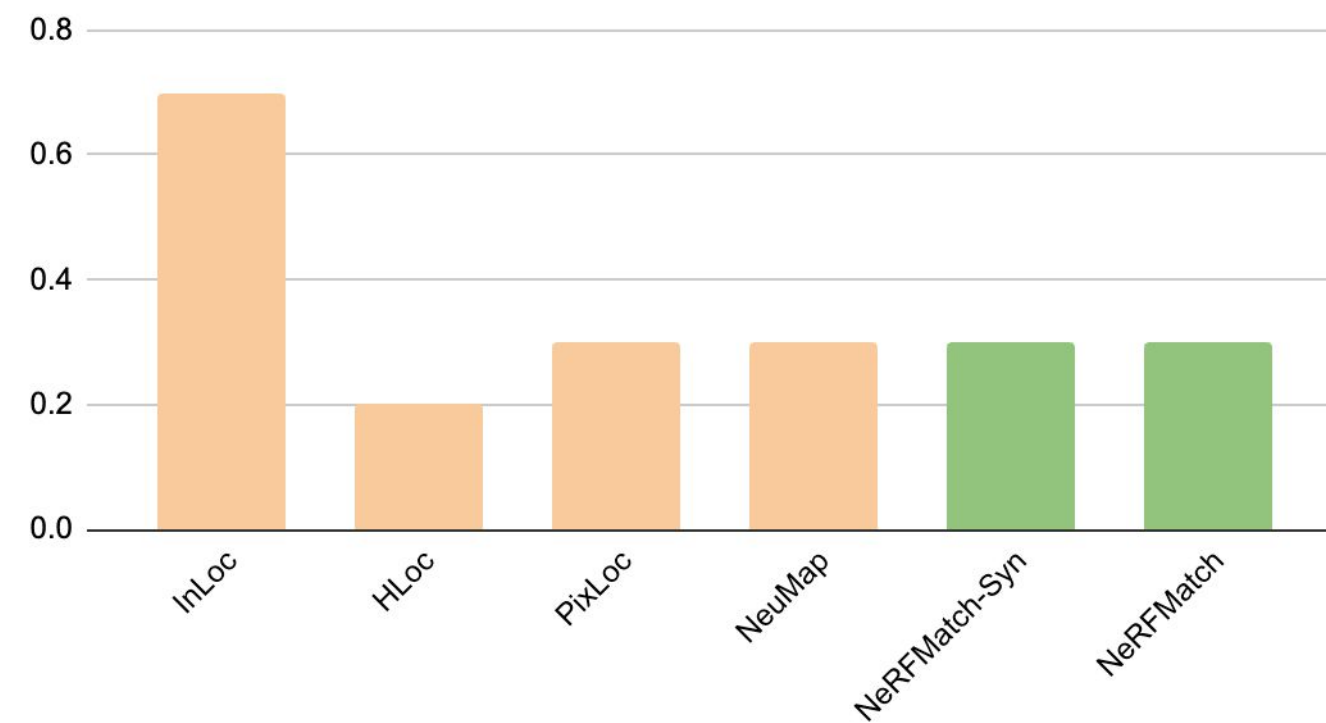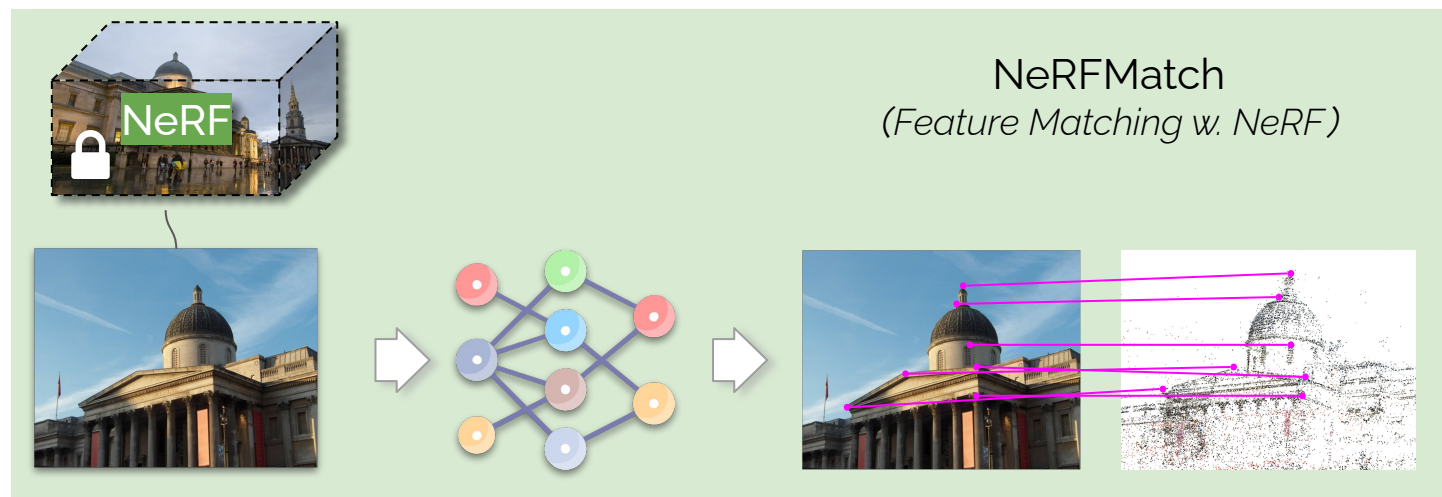
## Cambridge Landmarks - Translation Error (cm)



## Cambridge Landmarks - Rotation Error (°)

NeRFMatch
*(Feature Matching w. NeRF)*

VS

Feature Matching
*(w. 3D point cloud & images)*

**Cambridge Landmarks - Translation Error (cm)**

| | InLoc | HLoc | PixLoc | NeuMap | NeRFMatch-Syn | NeRFMatch |
|---|---|---|---|---|---|---|
| Value | ~49 | ~11 | ~15 | ~12 | ~11 | ~10.5 |

**Cambridge Landmarks - Rotation Error (°)**

| | InLoc | HLoc | PixLoc | NeuMap | NeRFMatch-Syn | NeRFMatch |
|---|---|---|---|---|---|---|
| Value | ~0.70 | ~0.20 | ~0.30 | ~0.30 | ~0.30 | ~0.30 |

- We are on-par with the SOTA HLoc results on Cambridge Landmarks, which is quite challenging wild environment for NeRF training.

NeRFMatch
*(Feature Matching w. NeRF)*

## Indoor performance bottleneck vs SOTA

- **Depth inaccuracies**: NeRF predicted depth maps are used to compute pseudo ground-truth for matching supervision. Incorrect depth predictions can lead to misaligned feature correspondences. In contrast, image matching, SCR, and APR methods use more accurate labels like Colmap camera poses or 3D maps.
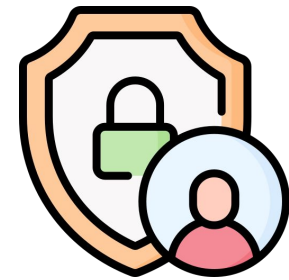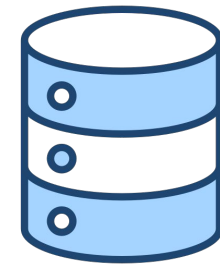
## Indoor performance bottleneck vs SOTA

- Not good yet at **filtering inaccurate matches**, which has a large effect on small scenes.

- Better **scaling** to large-outdoor scene compared to regression-based methods.

| Method | Scene Repres. | 7-Scenes - SfM Poses - Indoor | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Chess | Fire | Heads | Office | Pump. | Kitchen | Stairs | Avg.Med↓ | Avg.Recall↑. |
| MS-Trans. [53] | APR Net. | 11/6.4 | 23/11.5 | 13/13 | 18/8.1 | 17/8.4 | 16/8.9 | 29/10.3 | 18.1/9.5 | - |
| DFNet [17] | APR Net. | 3/1.1 | 6/2.3 | 4/2.3 | 6/1.5 | 7/1.9 | 7/1.7 | 12/2.6 | 6.4/1.9 | - |
| NeFeS [16] | APR+NeRF | 2/0.8 | 2/0.8 | 2/1.4 | 2/0.6 | 2/0.6 | 2/0.6 | 5/1.3 | 2.4/0.9 | - |
| DSAC* [10] | SCR Net. | 0.5/0.2 | 0.8/0.3 | 0.5/0.3 | 1.2/0.3 | 1.2/0.3 | 0.7/0.2 | 2.7/0.8 | 1.1/0.3 | **97.8** |
| ACE [6] | SCR Net. | 0.7/0.5 | 0.6/0.9 | 0.5/ 0.5 | 1.2/0.5 | 1.1/0.2 | 0.9/0.5 | 2.8/1.0 | 1.1/0.6 | 97.1 |
| DVLAD+R2D2 [45, 60] | 3D+RGB | **0.4/0.1** | **0.5/0.2** | **0.4/0.2** | **0.7/0.2** | **0.6/0.1** | **0.4/0.1** | **2.4/0.7** | **0.8/0.2** | 95.7 |
| HLoc [48] | 3D+RGB | 0.8/**0.1** | 0.9/**0.2** | 0.6/0.3 | 1.2/**0.2** | 1.4/0.2 | 1.1/**0.1** | 2.9/0.8 | 1.3/0.3 | 95.7 |
| NeRFMatch-Mini | NeRF+RGB | 1.4/0.5 | 1.7/1.0 | 2.1/0.7 | 4.4/1.0 | 4.7/1.0 | 2.2/0.5 | 8.8/2.1 | 3.6/0.9 | 67.9 |
| NeRFMatch | NeRF+RGB | 0.9/0.3 | 1.3/0.4 | 1.6/1.0 | 3.2/0.7 | 3.3/0.7 | 1.3/0.3 | 7.5/1.3 | 2.7/0.7 | 75.3 |
| NeRFMatch | NeRF | 0.9/0.3 | 1.3/0.4 | 1.6/1.0 | 3.3/0.7 | 3.2/0.6 | 1.3/0.3 | 7.2/1.3 | 2.7/0.7 | 75.4 |

# Conclusions

- Geometric localization is possible and (somewhat) SOTA

- Initial steps towards NeRF as the primary representation for visual localization
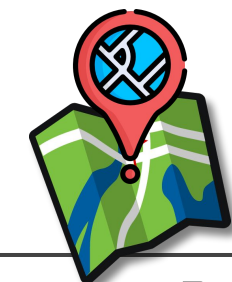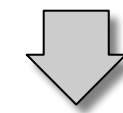
  ☺ Compact 🚀

  ☺ Interpretable
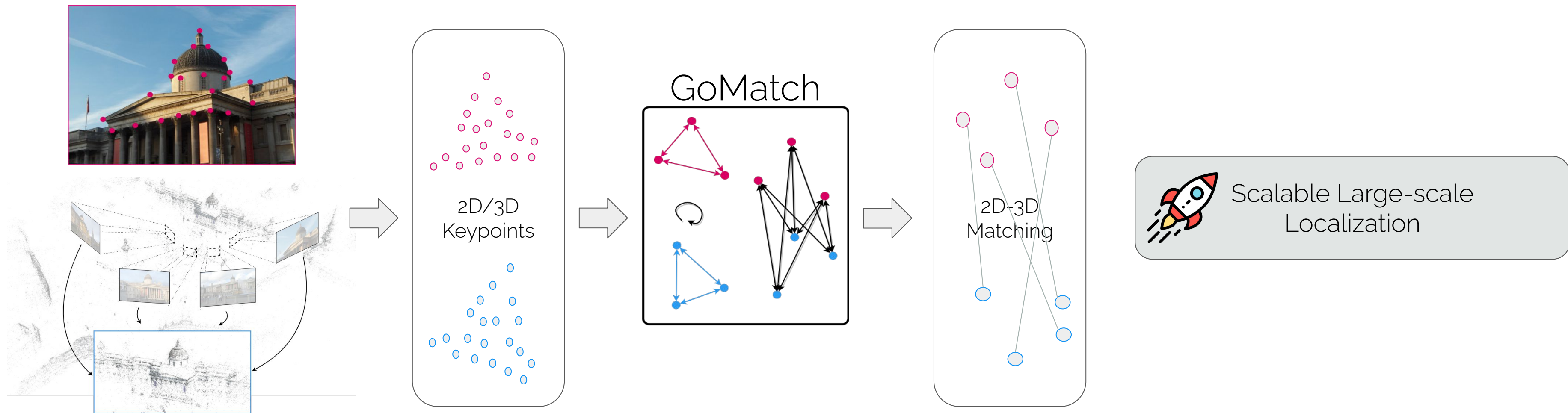
  ☺ 2D-3D matching

Camera Pose

Orientation
R

Position
(x, y, z)

Questions?

https://research.nvidia.com/labs/dvl/

Laura Leal-Taixé | CVPR | June 2024

# Things that did not make the cut (yet?)

GoMatch

2D/3D Keypoints

2D-3D Matching

Scalable Large-scale Localization

Qunjie Zhou

Sérgio Agostinho

Aljoša Ošep

Laura Leal-Taixé

# Visual Localization

# Visual Localization

# Visual Localization



Localize

# Overview



Low Storage

Privacy Preserving

No Descriptor Maintenance

Scalable Large-scale Localization

2D/3D Keypoints

GoMatch

2D-3D Matching

Localize

# Visual Localization



2D/3D Keypoints → Visual Descriptors → 2D-3D Matching

☺ Low Storage

☺ Privacy Preserving

☺ No Descriptor Maintenance

🚀 Scalable Large-scale Localization

# Practical Challenges

Storage Demand ⟶

| MegaDepth (192 scenes) | Storage | Desc Type | Data Type | Storage |
|---|---|---|---|---|
| Cameras | 15.73 MB | SIFT | Uint8 | 133.33 GB |
| 3D Points | 3.44 GB | CAPS | FP32 | 523.83 GB |
| Images | 157.84 GB | SuperPoint | FP32 | 1.044 TB |

# Practical Challenges

Storage Demand

| MegaDepth (192 scenes) | Storage | Desc Type | Data Type | Storage |
|---|---|---|---|---|
| Cameras | 15.73 MB | SIFT | Uint8 | 133.33 GB |
| 3D Points | 3.44 GB | CAPS | FP32 | 523.83 GB |
| Images | 157.84 GB | SuperPoint | FP32 | 1.044 TB |

3.44 GB

Scene Compression [1]

Desc Quantization [1, 2]

Scene Descriptor

133 GB ~ 1.04 TB

[1] Camposeco, Federico, et al. "Hybrid scene compression for visual localization." CVPR19
[2] Sattler, Torsten, Bastian Leibe, and Leif Kobbelt. "Efficient & effective prioritized matching for large-scale image-based localization." PAMI16

# Compare to E2E – Cambridge Landmarks



**Median Translation Error (m)** — GM / E2E / VM

HLoc 0.18, ActiveSearch 0.29, HybridSC 0.56, DSAC++ 0.14, MS-Transformer 1.28, MS-PoseNet 2.74, PoseNet 2.09, GoMatch 1.73, BPnPNet 11.44

**Median Rotation Error (deg)** — GM / E2E / VM

HLoc 0.63, ActiveSearch 0.63, HybridSC 0.66, DSAC++ 0.33, MS-Transformer 2.73, MS-PoseNet 5.34, PoseNet 6.84, GoMatch 11.02, BPnPNet 106.72

**Storage (MB)** — GM / E2E / VM

HLoc 3214.84, ActiveSearch 812.7, HybridSC 3.13, DSAC++ 828, MS-Transformer 71.1, PoseNet 200, GoMatch 48.15, BPnPNet 48.15

# Existing Solutions

Storage / Memory Efficiency

Privacy Preserving

GoMatch

Descriptor Maintenance

# Geometric-based matching and pose estimation

**SoftPOSIT [1]**
- Alternate step: softassign + POSIT
- Requires initialization
- Struggles with clutter, occlusions, repetitive patterns.
- Efficient



**GOPAC [3]**
- Globally optimal solution using Branch-and-Bound
- Prohibitive runtime requirements
- Cannot scale to large problems



hegyhati.github.io

**Bind PnP [2]**
- Kalman-Filter to maintain correspondence hypotheses.
- Requires initialization of GMM pose priors
- Better handling of occlusion, clutter and repetitive patterns



**BPnPNet [4]**
- Learning-based geometric matching network
- Declarative layers to back propagate through Sinkhorn, RANSAC and the PnP solver.
- Performance substantially degraded with outliers.

[1] David, Philip, et al. "SoftPOSIT: Simultaneous pose and correspondence determination." IJCV 2004
[2] Moreno-Noguer, Francesc et al. "Pose priors for simultaneously solving alignment and correspondence." ECCV 2008
[3] Campbell, Dylan, et al. "Globally-optimal inlier set maximisation for camera pose and correspondence estimation." PAMI 2018
[4] Campbell, Dylan, et al. "Solving the blind perspective-n-point problem end-to-end with robust differentiable geometric optimization." ECCV 2020.

# Geometric-Only Methods



Liu Liu, et al. "Learning 2d-3d correspondences to solve the blind perspective-n-point problem." arXiv20

Dylan Campbell, Liu Liu, and Stephen Gould. "Solving the blind perspective-n-point problem end-to-end with robust differentiable geometric optimization." ECCV20
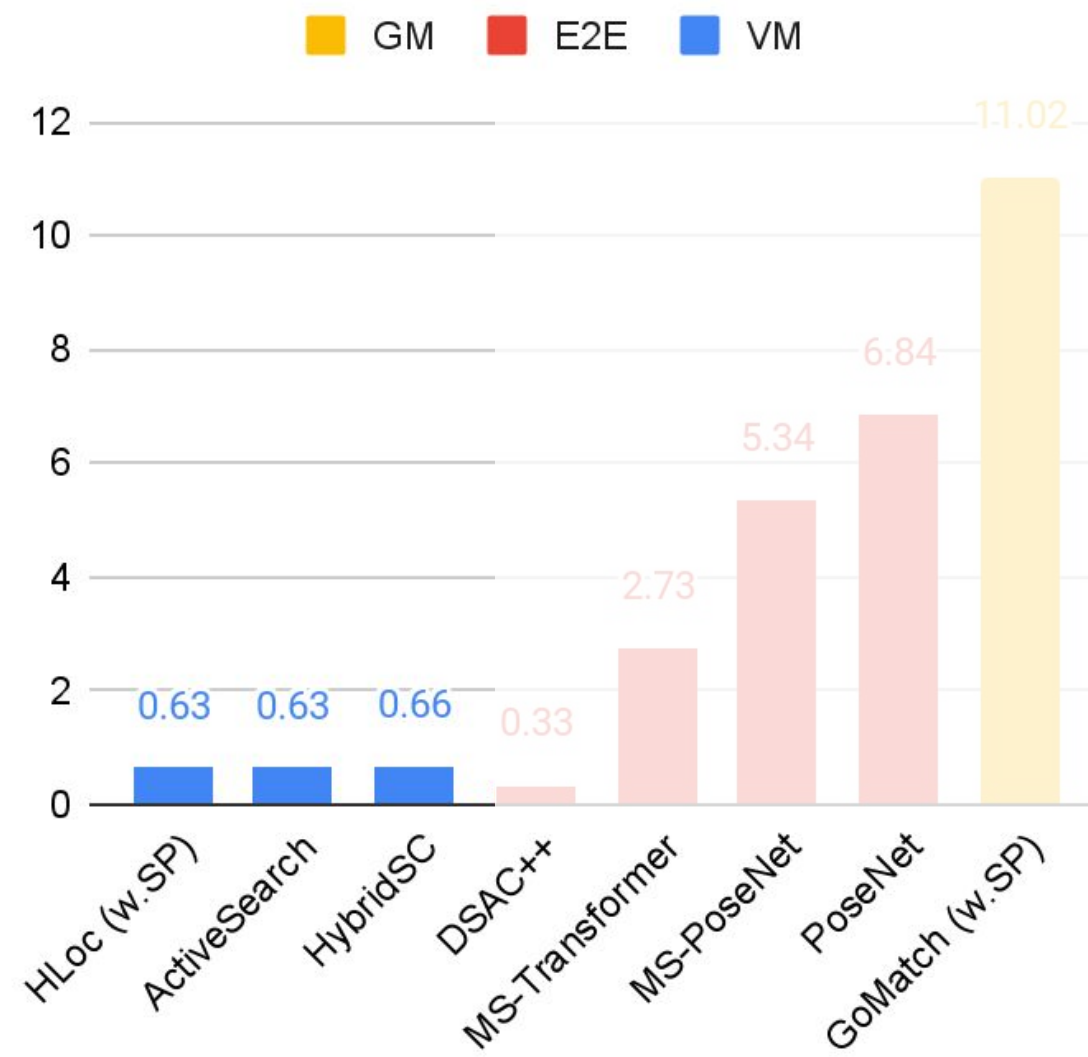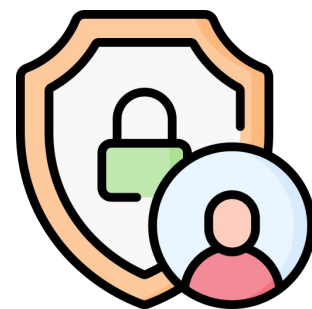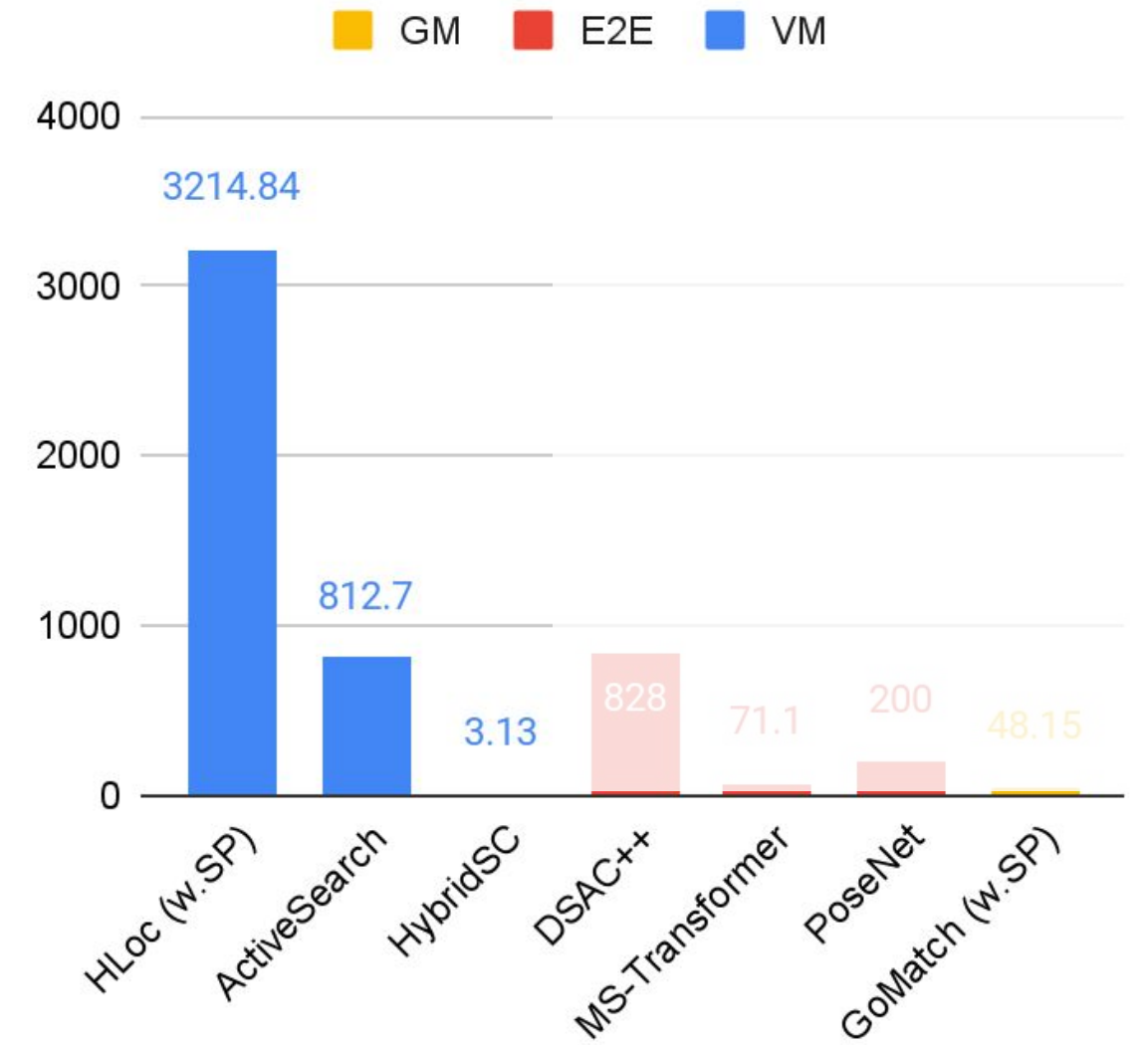
# Outdoor Scene – Cambridge Landmarks



**Median Translation Error (m)** — GM, E2E, VM

| Method | Value |
|---|---|
| HLoc (w.SP) | 0.18 |
| ActiveSearch | 0.29 |
| HybridSC | 0.56 |
| DSAC++ | 0.14 |
| MS-Transformer | 1.28 |
| MS-PoseNet | 2.74 |
| PoseNet | 2.09 |
| GoMatch (w.SP) | 1.73 |

**Median Rotation Error (deg)** — GM, E2E, VM

| Method | Value |
|---|---|
| HLoc (w.SP) | 0.63 |
| ActiveSearch | 0.63 |
| HybridSC | 0.66 |
| DSAC++ | 0.33 |
| MS-Transformer | 2.73 |
| MS-PoseNet | 5.34 |
| PoseNet | 6.84 |
| GoMatch (w.SP) | 11.02 |

**Storage (MB)** — GM, E2E, VM

| Method | Value |
|---|---|
| HLoc (w.SP) | 3214.84 |
| ActiveSearch | 812.7 |
| HybridSC | 3.13 |
| DSAC++ | 828 |
| MS-Transformer | 71.1 |
| PoseNet | 200 |
| GoMatch (w.SP) | 48.15 |

GoMatch

# Practical Challenges

Storage Demand

| MegaDepth | Storage | Desc Type | Data Type | Storage |
|---|---|---|---|---|
| **Cameras** | 15.73 MB | SIFT | Uint8 | 133.33 GB |
| **3D Points** *(192 scenes)* | 3.44 GB | CAPS | FP32 | 523.83 GB |
| **Images** *(192 scenes)* | 157.84 GB | SuperPoint | FP32 | 1.044 TB |

3.44 GB

Scene Compression [1]

Desc Quantization [1, 2]

Scene Descriptor

133 GB ~ 1.04 TB

[1] Camposeco, Federico, et al. "Hybrid scene compression for visual localization." CVPR19
[2] Sattler, Torsten, Bastian Leibe, and Leif Kobbelt. "Efficient & effective prioritized matching for large-scale image-based localization." PAMI16

GoMatch

# Practical Challenges



Privacy Risk

Descriptor Inversion [1]

Client

Matching + Pose Estimation

Server

Matching + Pose Estimation

Privacy-preserving Descriptors [2, 3]

Keypoints    Subspaces    Inversion    Reconstruction

NinjaNet

input image    original descriptors    content-concealing descriptors

[1] Francesco, Pittaluga, et al Revealing Scenes by Inverting Structure From Motion Reconstructions. CVPR19
[2] Dusmanu, Mihai, et al. "Privacy-preserving image features via adversarial affine subspace embeddings." CVPR21.
[3] Ng, Tony, et al. "NinjaDesc: Content-Concealing Visual Descriptors via Adversarial Learning." CVPR22

# Practical Challenges



Maintenance Complexity

Descriptor Upgrade [1]

SIFT → CAPS → SuperPoint

Cross-device Matching and Localization [1]

SIFT

CAPS

HardNet

SOSNet

SuperPoint

[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Visual Localization



Localize

# GoMatch Step-by-Step

# Introduction



Localize

2D/3D Keypoints

Visual Descriptors

2D-3D Matching

# Overview



Localize

2D/3D
Keypoints

Visual
Descriptors

2D-3D
Matching

# Overview



Localize

2D/3D Keypoints

GoMatch

2D-3D Matching

🙂 Low Storage

🙂 Privacy Preserving

🙂 No Descriptor Maintenance

79

# Overview



Localize



2D/3D Keypoints → GoMatch → 2D-3D Matching

Low Storage

Privacy Preserving

No Descriptor Maintenance

Scalable Large-scale Localization
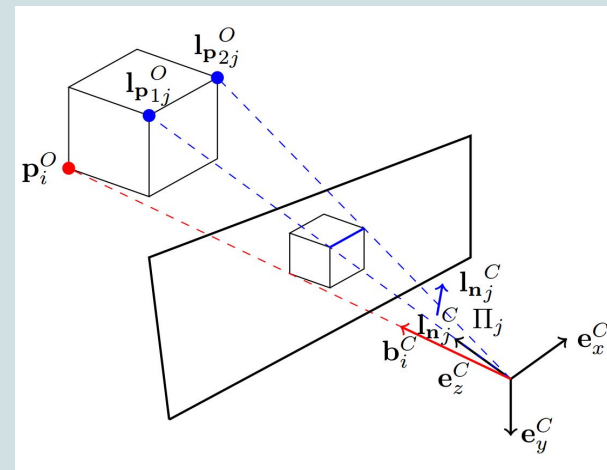
# Motivation

## End-to-end



Camera Pose Regression
(PoseNet, MapNet, …)

Scene Coordinate Regression
(DSAC, DSAC++, …)

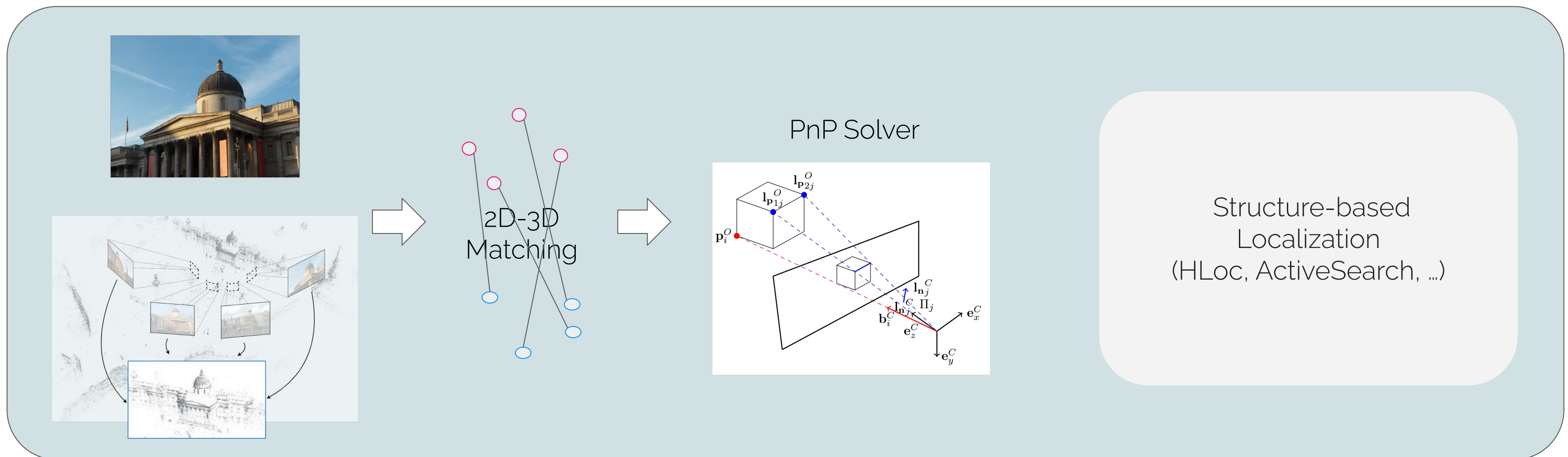## Structure-based



2D-3D Matching
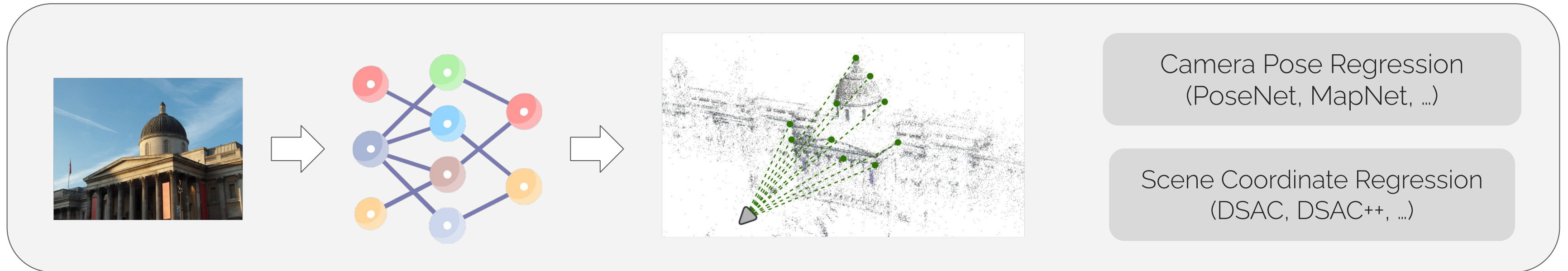
Perspective-n-Point Solver

Structure-based Localization
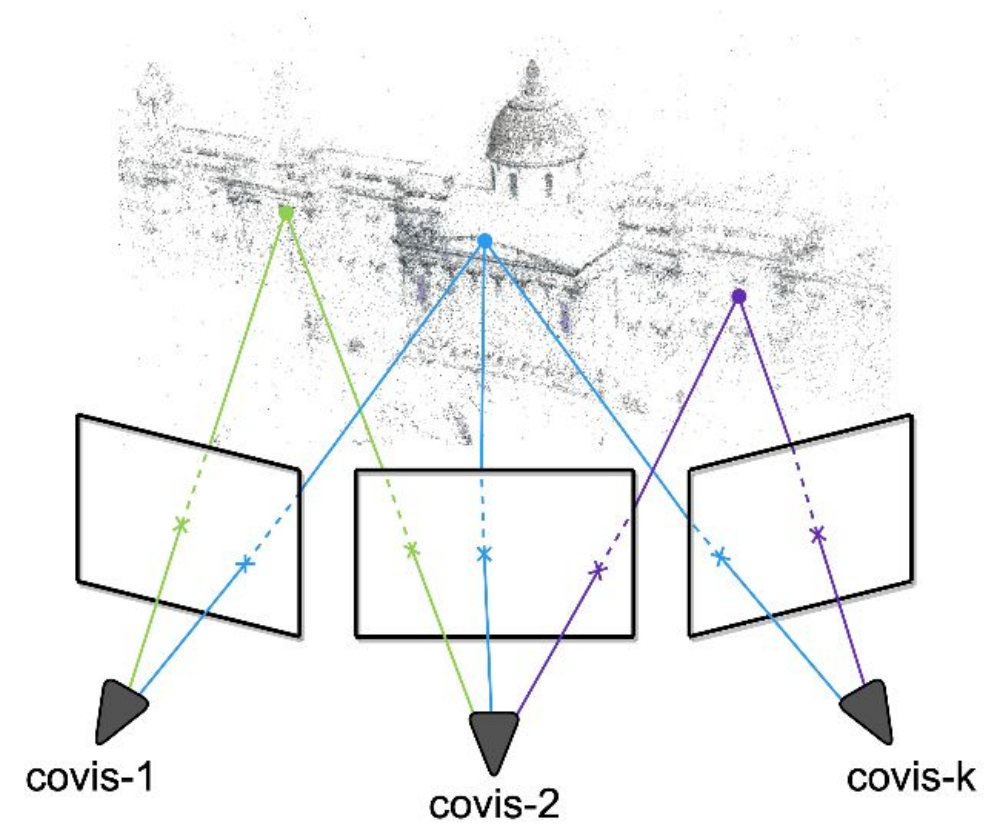(HLoc, ActiveSearch, …)

# Visual Localization Approaches



Camera Pose Regression
(PoseNet, MapNet, …)

Scene Coordinate Regression
(DSAC, DSAC++, …)

2D-3D Matching

PnP Solver

Structure-based Localization
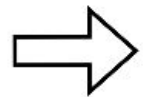(HLoc, ActiveSearch, …)

# GoMatch



Query Image
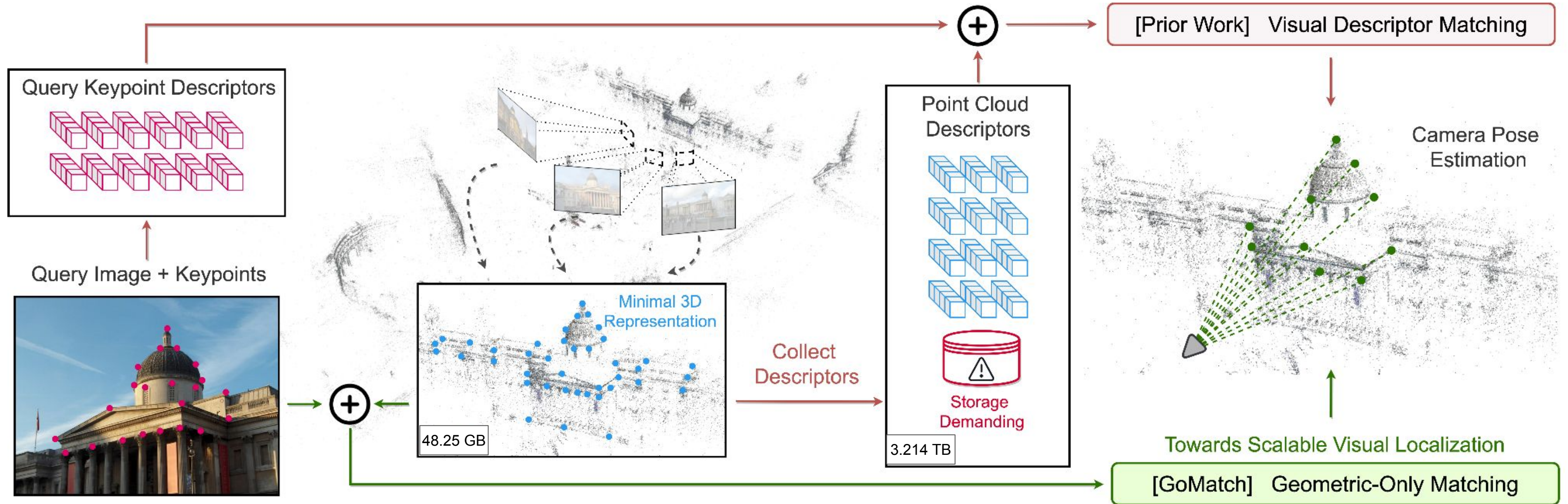
Image Retrieval

covis-1  covis-2  covis-k

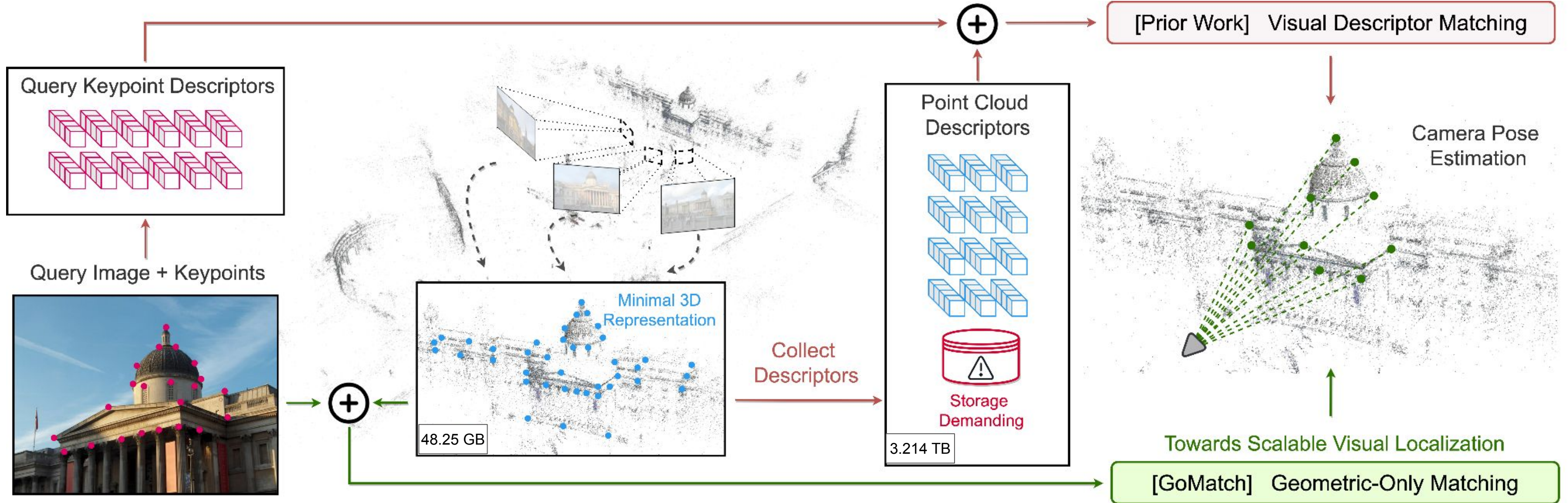Localization Performance

Storage Requirements

Privacy

No Descriptor Maintenance

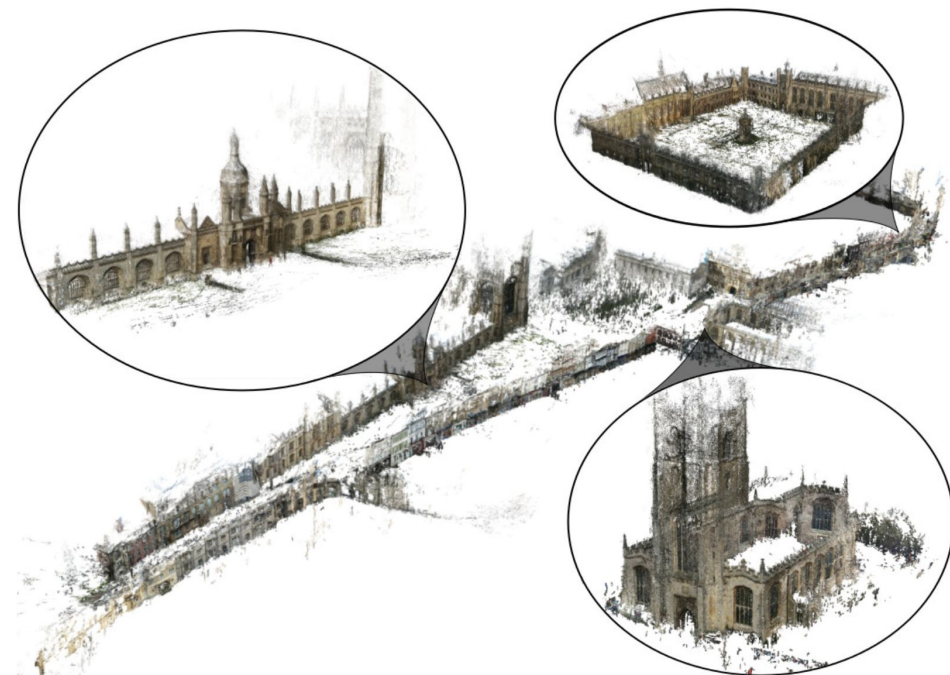# Significantly Lower Storage Requirements



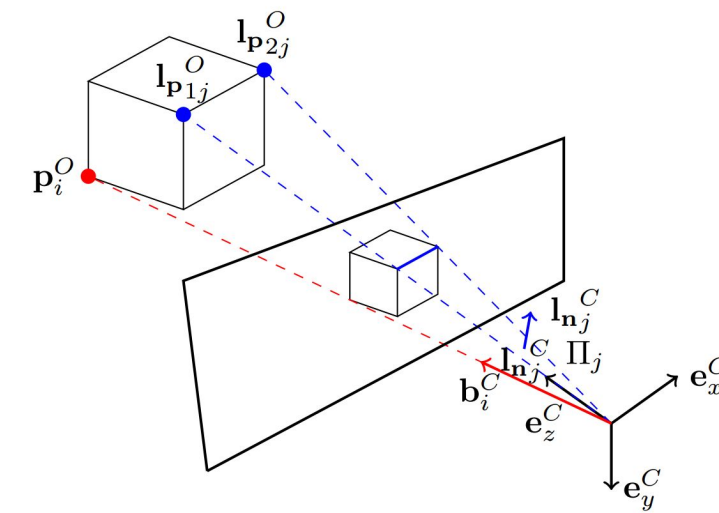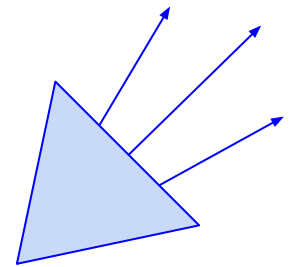1.5% vs visual descriptors

# Classical Structure-based Localization

# Structure-based Approaches



Scene representation

2D-3D Correspondences

PnP Solver

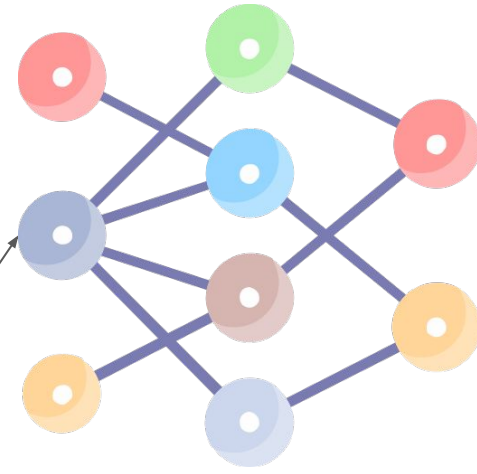T. Sattler, B. Leibe, and L. Kobbelt. Efficient & Effective Prioritized Matching for Large-Scale Image-Based Localization. PAMI2017

# End-to-end Learned Localization

Sattler, Torsten, Qunjie Zhou, Marc Pollefeys, and Laura Leal-Taixé. "Understanding the limitations of cnn-based absolute camera pose regression." CVPR19.

Scene Coordinate Regression
(DSAC, DSAC++, …)

Camera Pose Regression
(PoseNet, MapNet, …)

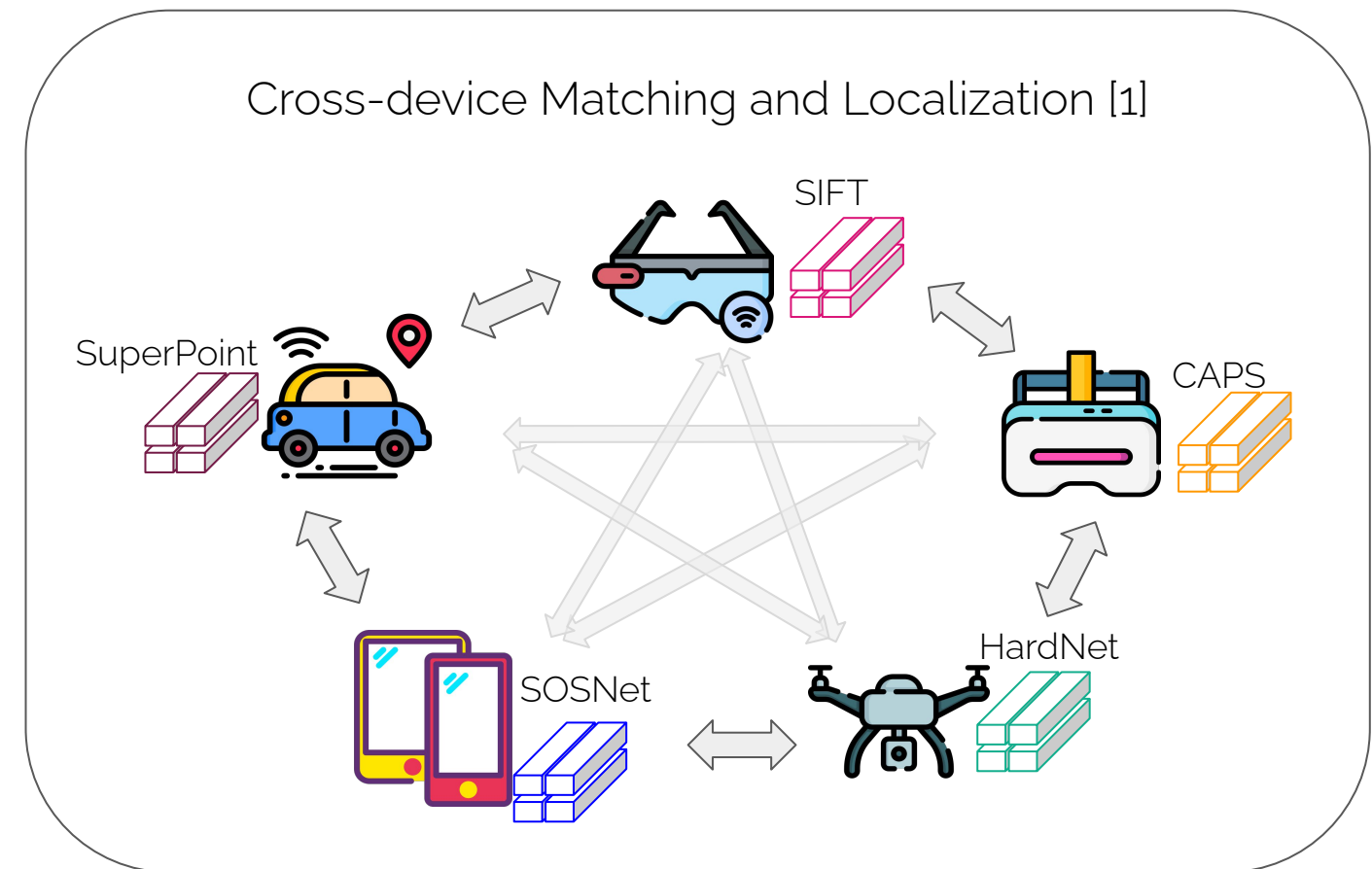Scene representation

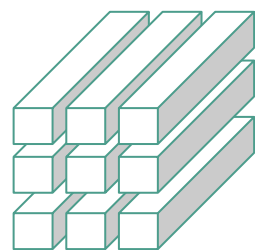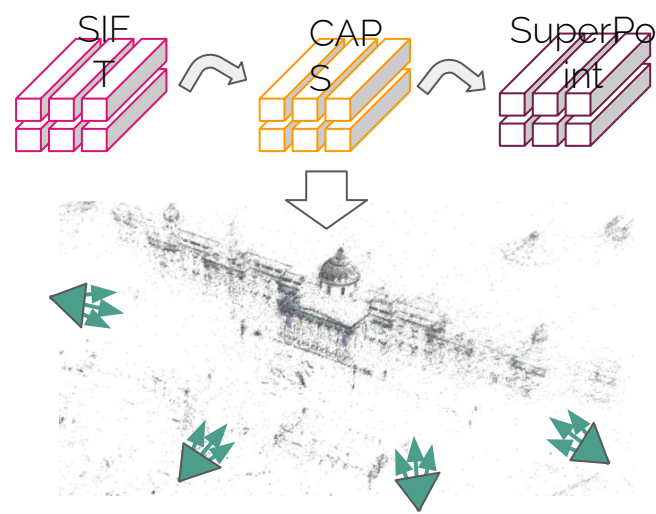Image Retrieval
(Netvlad, GeM, …)

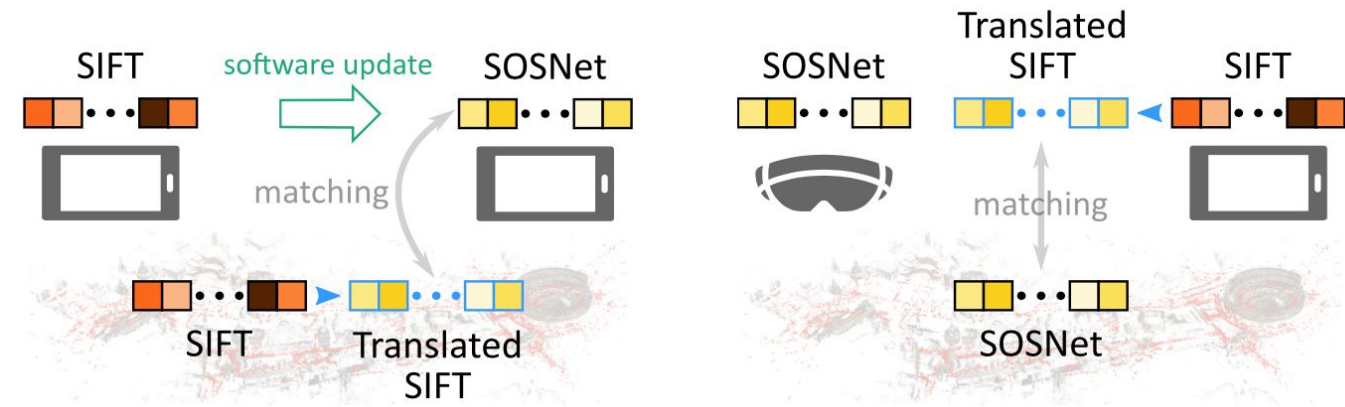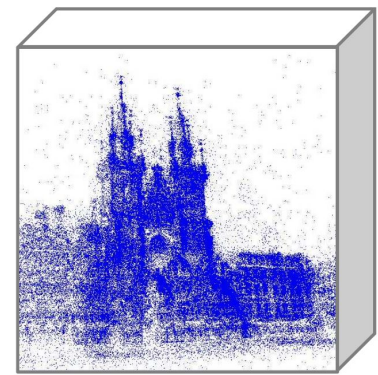Relative Pose Estimation
(EssNet, CamNet, …)

# Storage Requirements

# Practical Challenges

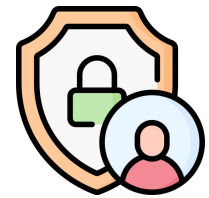Maintenance
Complexity

SIFT → software update → SOSNet
matching

SOSNet    Translated SIFT    SIFT
matching

SIFT → Translated SIFT

SOSNet

SIFT ⤳ CAPS ⤳ SuperPoint

Cross-device Matching and Localization [1]

SIFT

SuperPoint

CAPS

SOSNet

HardNet
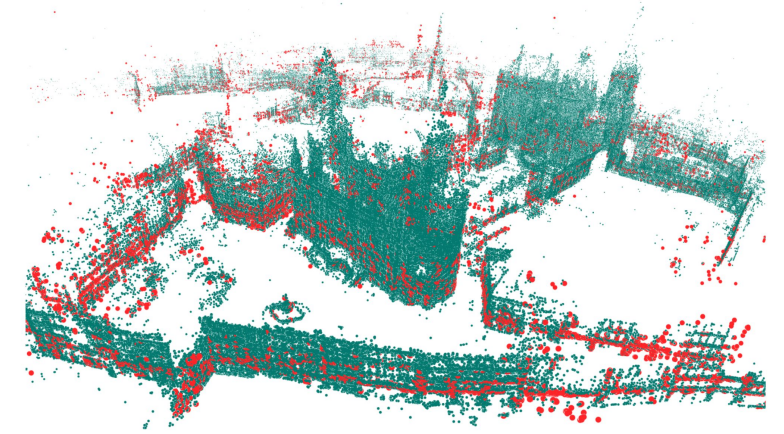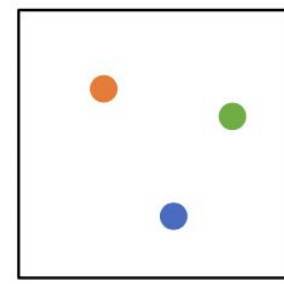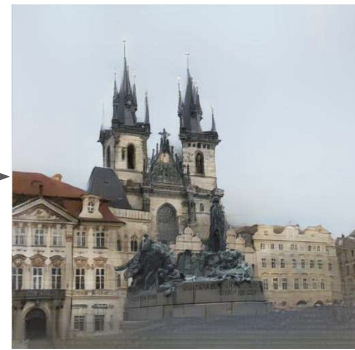
[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Practical Challenges



SIFT → matching → BRIEF

↕ localization and mapping ↕

SOSNet ← HardNet

Descriptor Inversion

Keypoints · Subspaces → Inversion → Reconstruction · input image

NinjaNet

original descriptors · content-concealing descriptors

SIFT → software update → SOSNet

matching

SIFT → Translated SIFT

SOSNet · Translated SIFT · SIFT

matching

SOSNet

# Practical Challenges



Scene Descriptor

Storage Demand

Privacy Risk

Client

Server

Maintenance Complexity

SIFT

CAPS

SuperPoint

SuperPoint

SIFT

CAPS

HardNet

SOSNet

# Existing Solutions



GoMatch

Storage / Memory Efficiency

Descriptor Maintenance

Privacy Preserving

# Practical Challenges



Maintenance Complexity

Descriptor Upgrade [1]

SIFT          CAPS          SuperPoint

Map Re-building

Upgraded Scene Descriptor

Cross-device Matching and Localization [1]

SIFT

SuperPoint

CAPS

SOSNet

HardNet

[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Practical Challenges



Maintenance Complexity

Descriptor Upgrade [1]

SIFT    CAPS    SuperPoint

Map Re-building

Upgraded Scene Descriptor

Cross-device Matching and Localization [1]

SIFT    software update    SOSNet

SIFT    Translated SIFT

SOSNet    Translated SIFT    SIFT

SOSNet

matching    matching

SuperPoint    SIFT    CAPS

SOSNet    HardNet

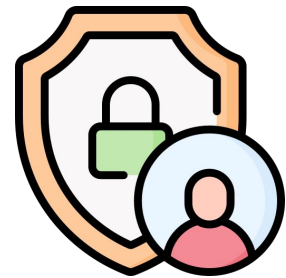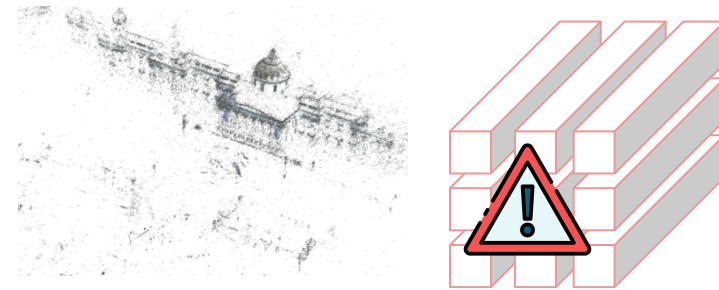[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Practical Challenges
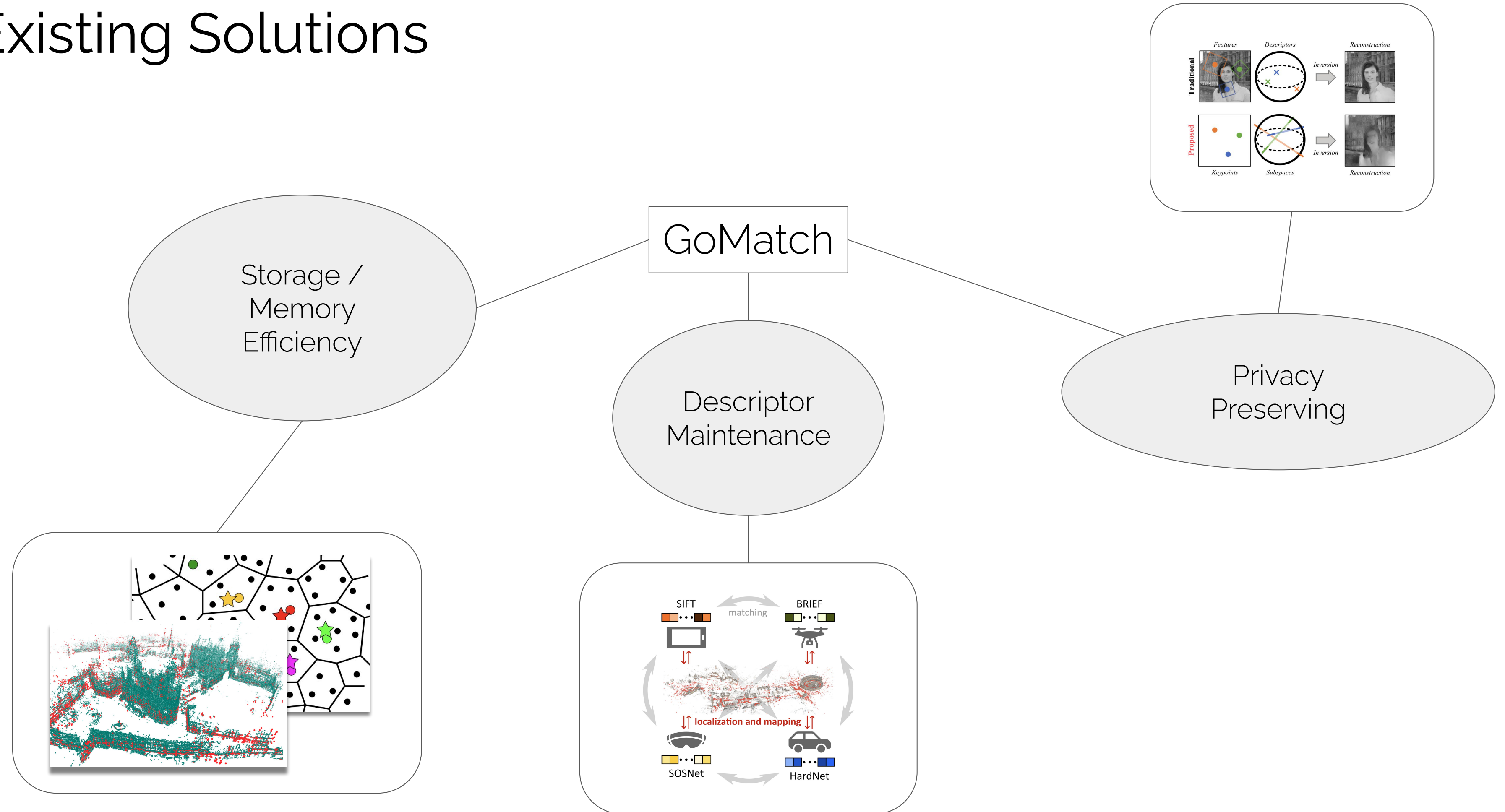


Storage Demand

Scene Descriptor

Scene Compression

Desc Quantization

Privacy Risk

Client

Server

Privacy-preserving Descriptors

Keypoints    Subspaces    original descriptors    content-concealing descriptors

Maintenance Complexity

SIFT    CAPS    SuperPoint

SuperPoint    SIFT    CAPS    HardNet    SOSNet

Descriptors Translation

SIFT    software update    SOSNet    SOSNet    Translated SIFT    SIFT

matching

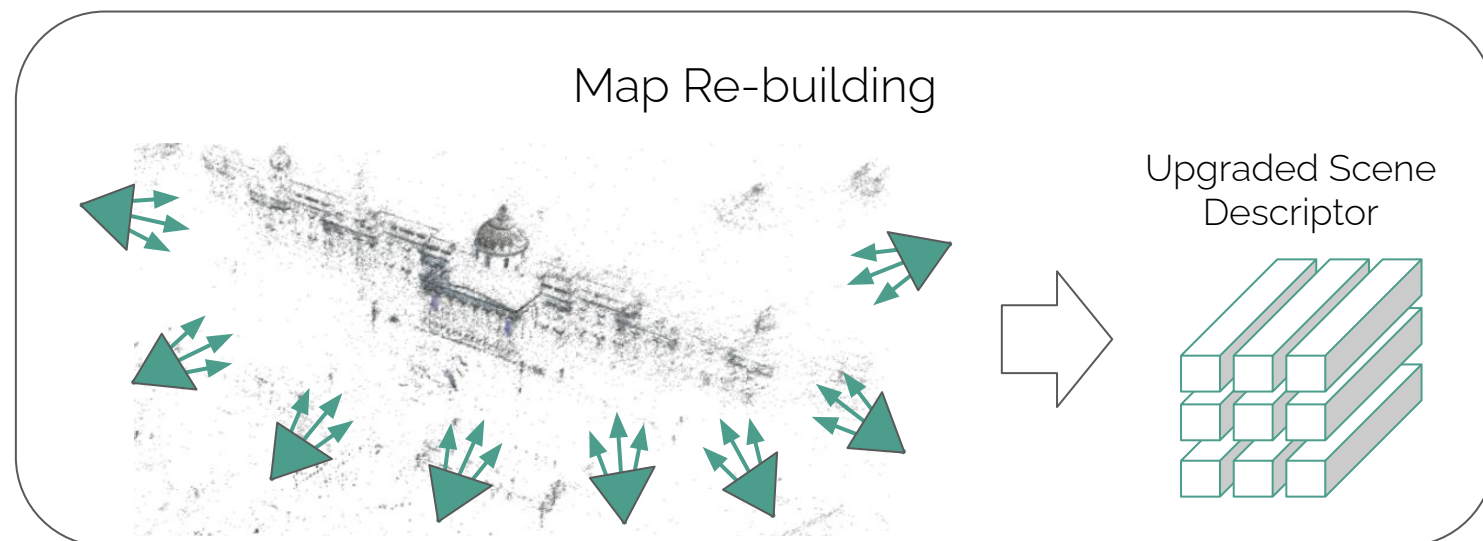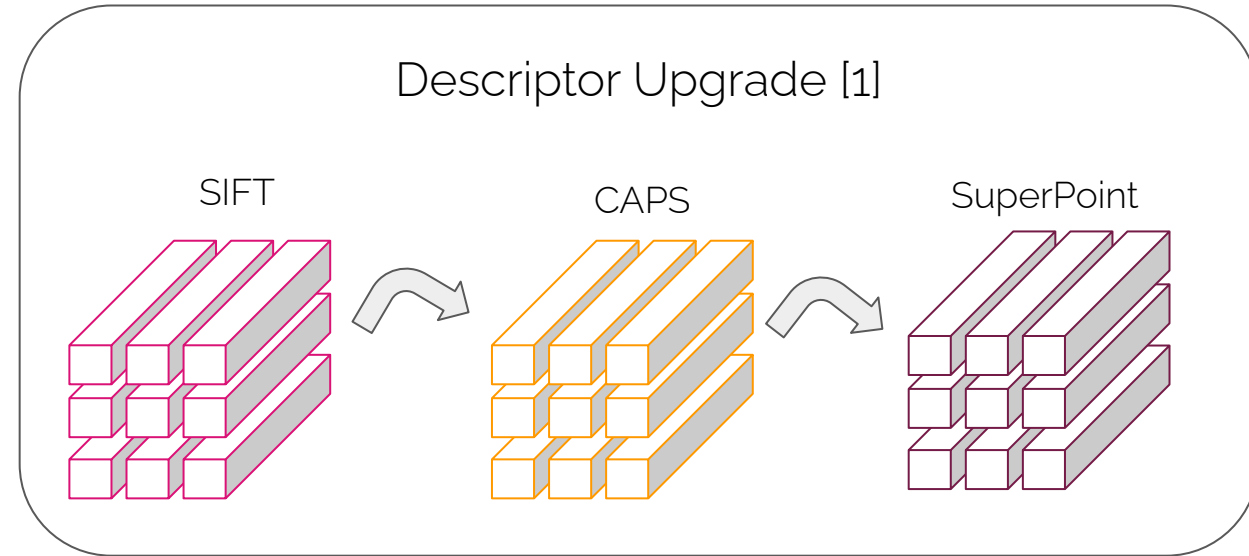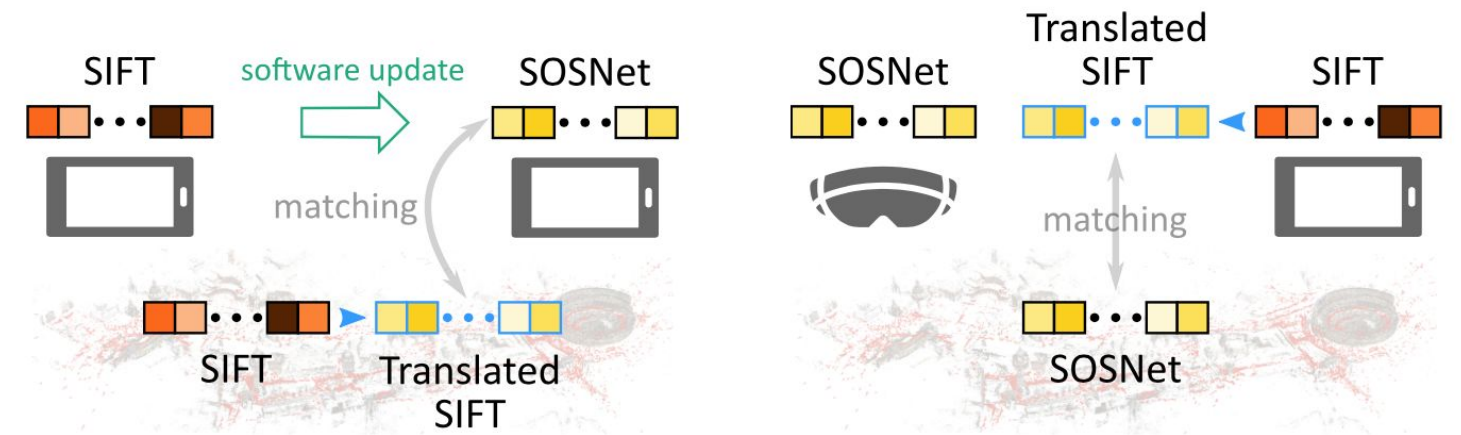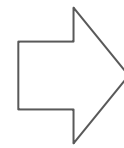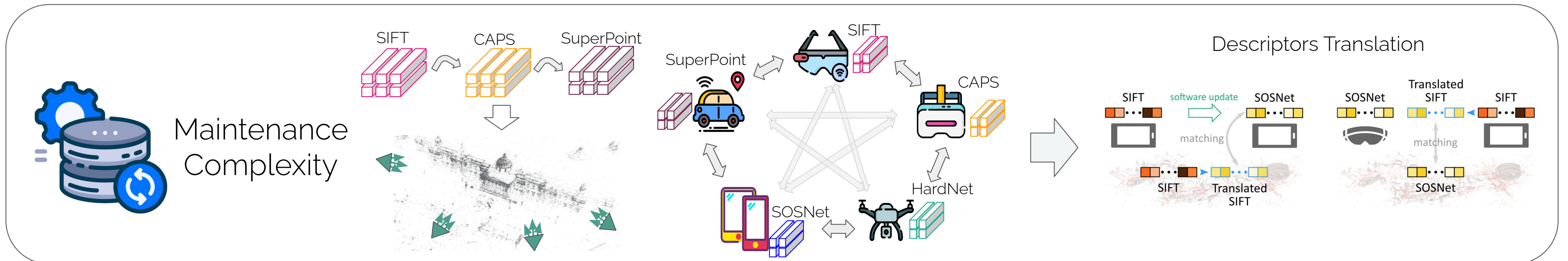SIFT    Translated SIFT    SOSNet

matching

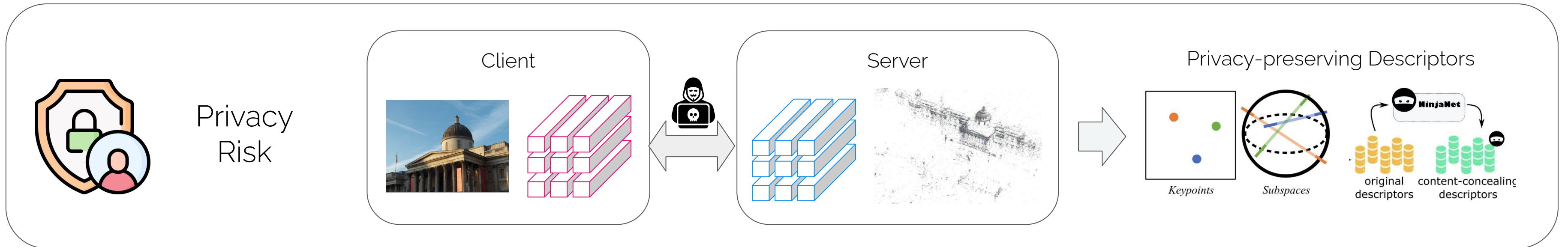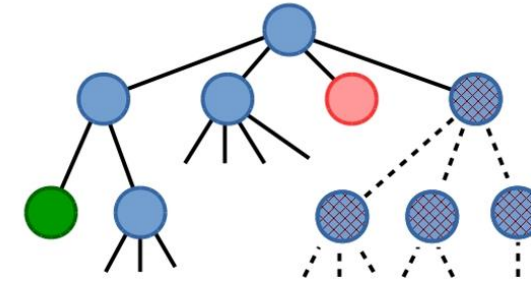# Maintenance Effort



[1] Dusmanu, Mihai, et al.. Cross-descriptor visual localization and mapping. ICCV21

# Descriptor Maintenance

# Geometric-based matching and pose estimation



hegyhati.github.io

[1] David, Philip, et al. "SoftPOSIT: Simultaneous pose and correspondence determination." IJCV 2004
[2] Moreno-Noguer, Francesc et al. "Pose priors for simultaneously solving alignment and correspondence." ECCV 2008
[3] Campbell, Dylan, et al. "Globally-optimal inlier set maximisation for camera pose and correspondence estimation." PAMI 2018
[4] Campbell, Dylan, et al. "Solving the blind perspective-n-point problem end-to-end with robust differentiable geometric optimization." ECCV 2020.

SoftPOSIT [1]
- Alternate step: softassign + POSIT
- Requires initialization
- Struggles with clutter, occlusions, repetitive patterns.
- Efficient

GOPAC [3]
- Globally optimal solution using Branch-and-Bound
- Prohibitive runtime requirements
- Cannot scale to large problems

Bind PnP [2]
- Kalman-Filter to maintain correspondence hypotheses.
- Requires initialization of GMM pose priors
- Better handling of occlusion, clutter and repetitive patterns

BPnPNet [4]
- Learning-based geometric matching network
- Declarative layers to back propagate through Sinkhorn, RANSAC and the PnP solver.
- Performance substantially degraded in the presence of outliers.

# Geometric-based matching and pose estimation

SoftPOSIT [1]
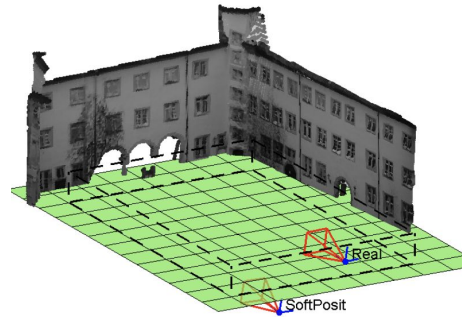- Iterative softassign
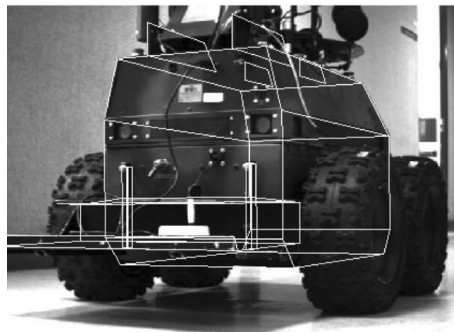- POSIT
- Requires initialization
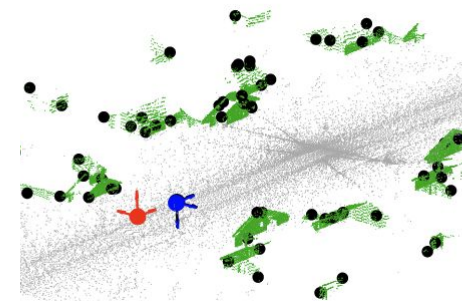
Bind PnP
Moreno-Noguer et al.
2008



BPnPNet [4]
Campbell et al. 2018



SoftPOSIT
Dementhon et al.
2004



GOPAC
Campbell et al. 2018



[1] David, Philip, et al. "SoftPOSIT: Simultaneous pose and correspondence determination." IJCV
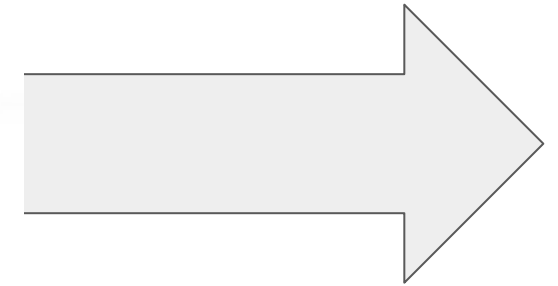[

# Geometric-based matching



Existing work

**With 2D–3D correspondences:**

- Perspective-n-Point (PnP)
  - Gao *et al.* 2003; Lepetit *et al.* 2009
  - + RANSAC [Fischler & Bolles 1981]
  - + global optimisation [Li 2009]
  - + neural network [Dang *et al.* 2018]

- Sparse feature pipelines
  - Svärm *et al.* 2016; Sattler *et al.* 2017; Cavallari *et al.* 2017, 2019; Schönberger *et al.* 2018; Taira *et al.* 2018
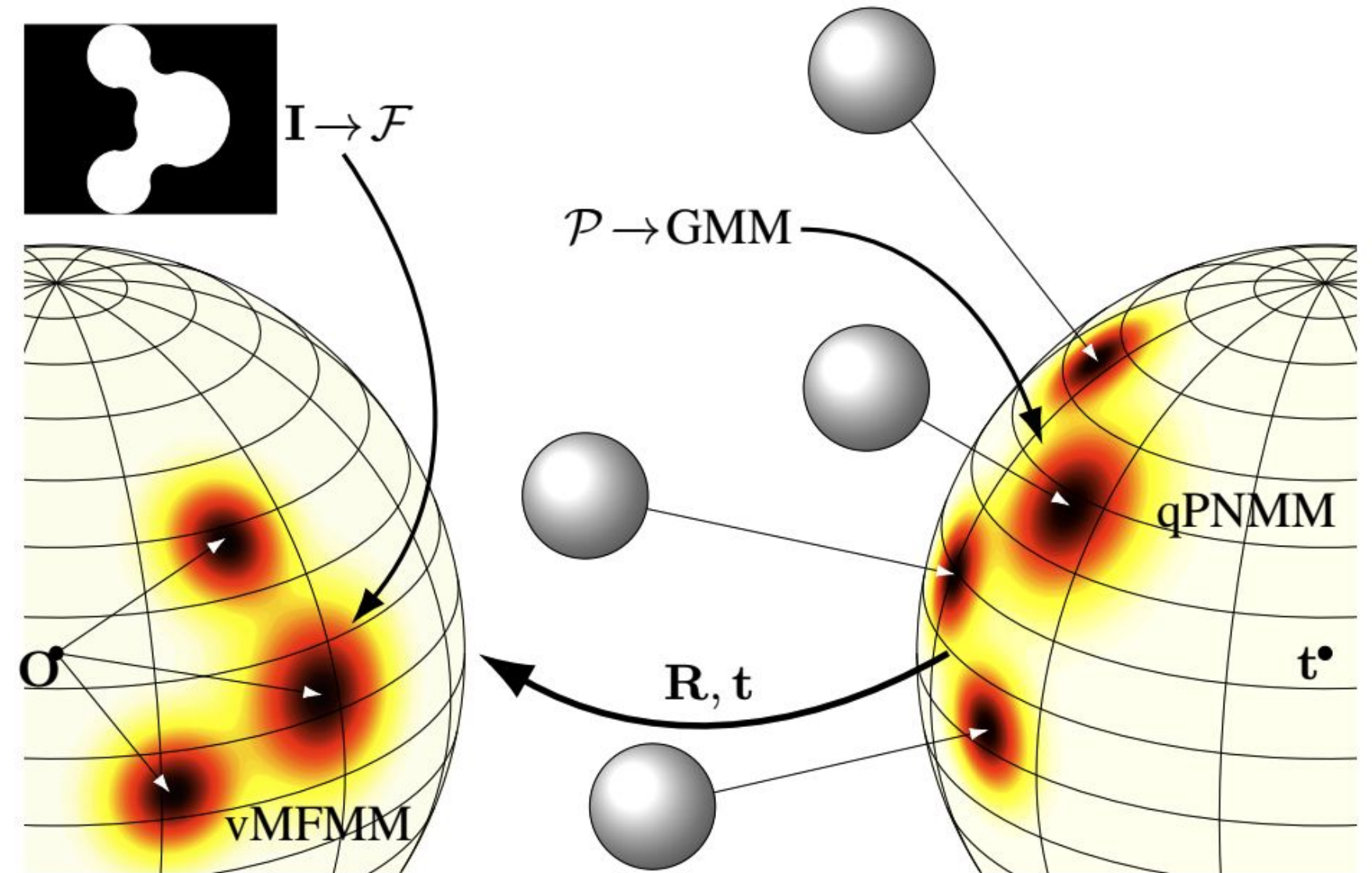
**Without 2D–3D correspondences:**

- Learning-based camera pose
  - Kendall *et al.* 2015–2017; Cai *et al.* 2018; Brahmbhatt *et al.* 2018; Radwan *et al.* 2018; Walch *et al.* 2017; Brachmann *et al.* 2017, 2018, 2020 (DSAC)

- Optimization-based camera pose
  - **Local:** David et al. 2004 (SoftPOSIT); Moreno-Noguer et al. 2008 (BlindPnP)
  - **Global:** Grimson 1990; Jurie 1999; Brown *et al.* 2015; Campbell *et al.* 2019

# Geometric-Only Methods



Campbell, Dylan, Lars Petersson, Laurent Kneip, Hongdong Li, and Stephen Gould. The alignment of the spheres: Globally-optimal spherical mixture alignment for camera pose estimation. CVPR 2019.
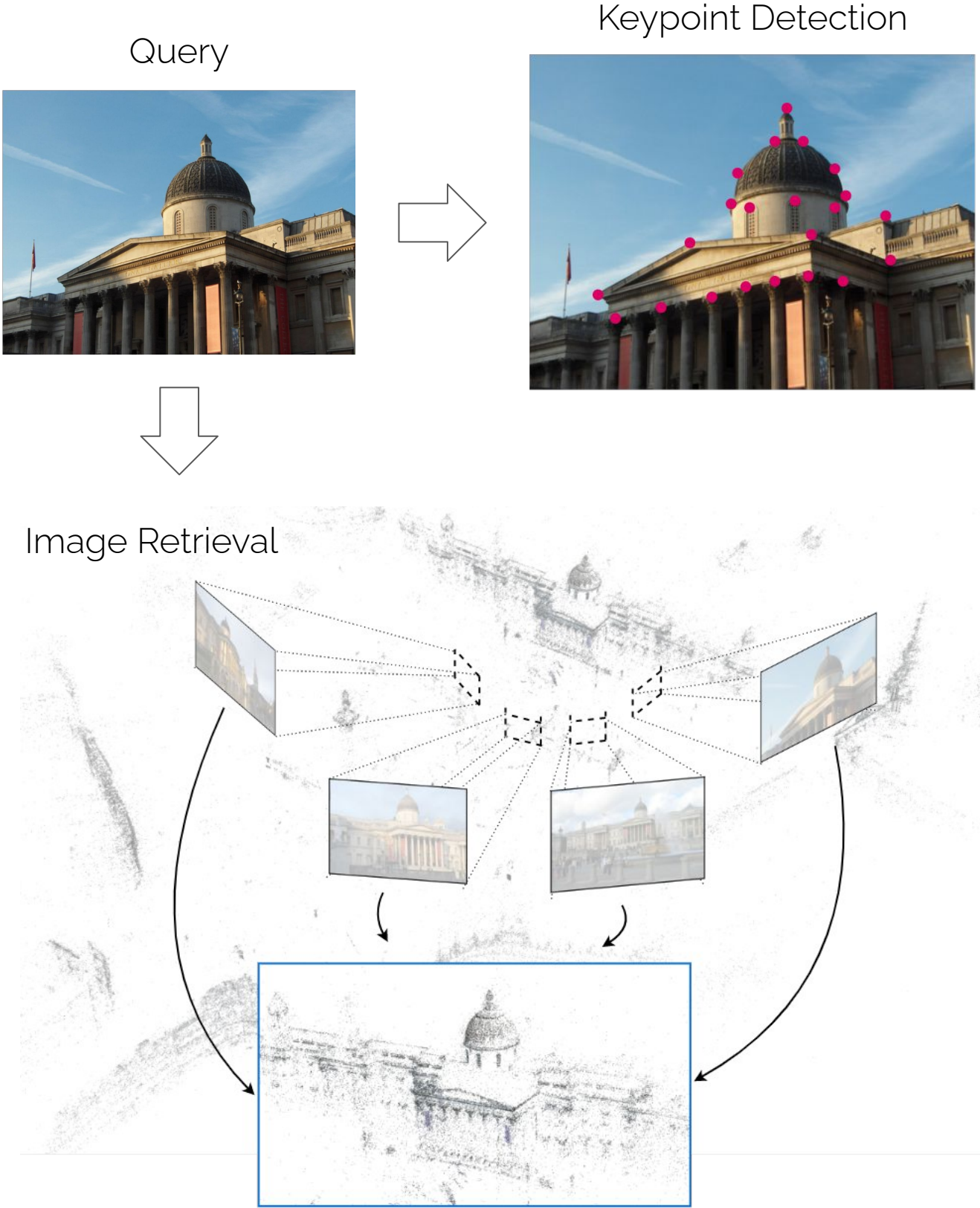
# Visual Localization

# Classical Structure-based Localization

Query

# Classical Structure-based Localization



Query

Keypoint Detection

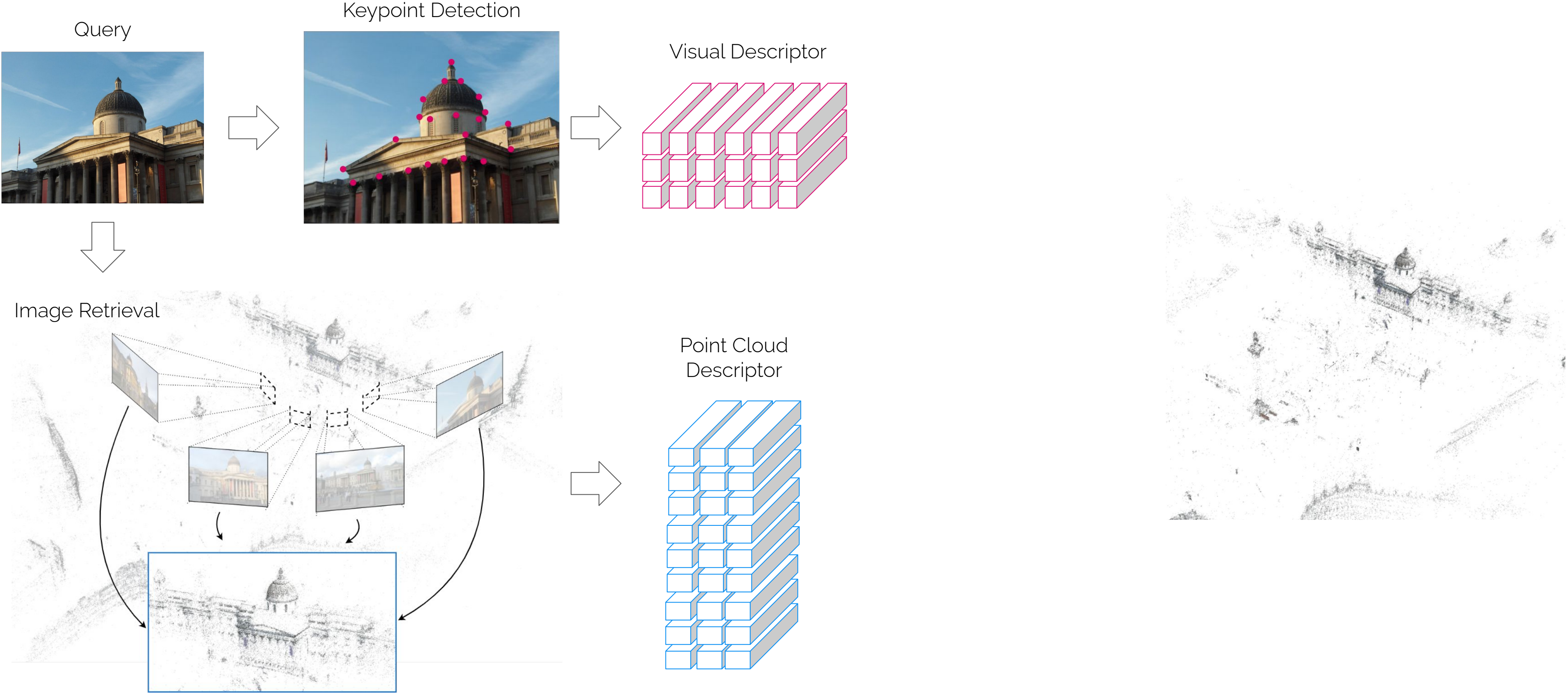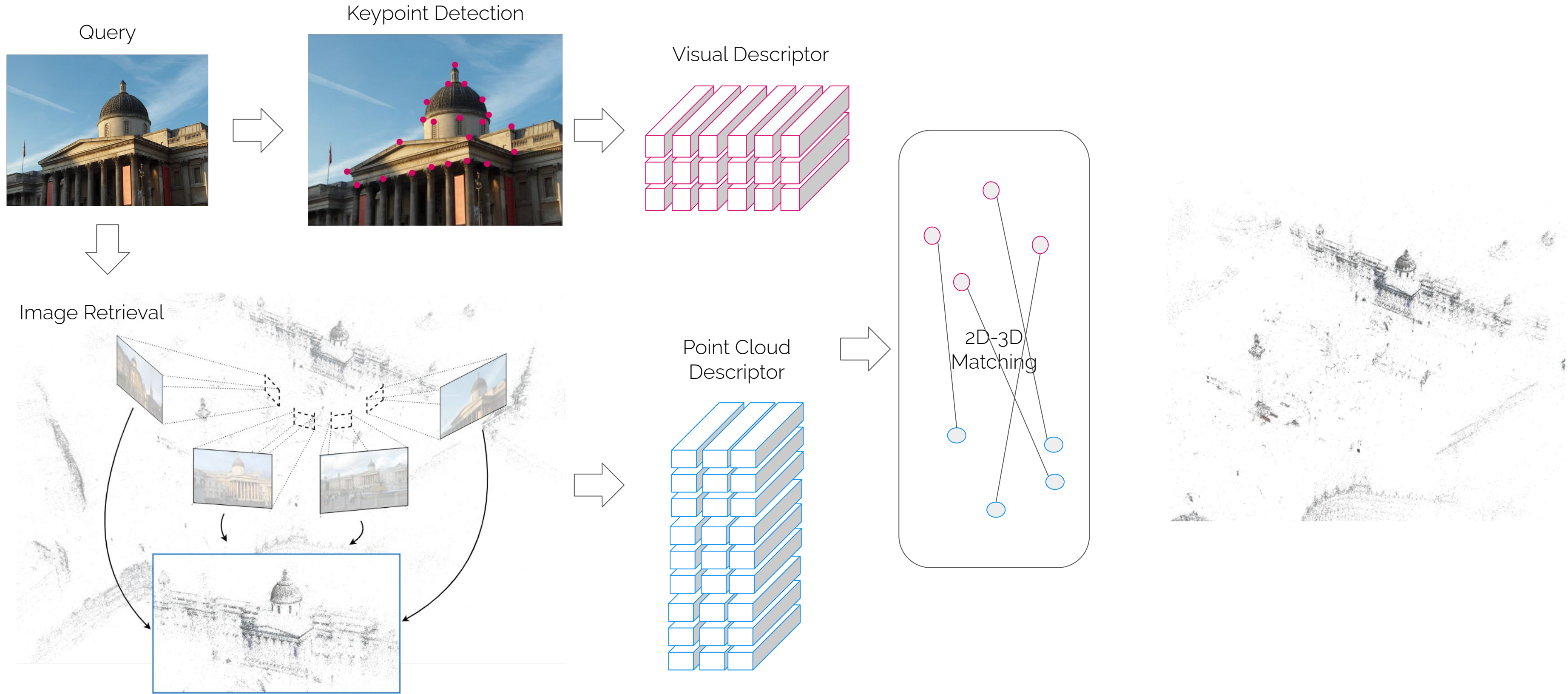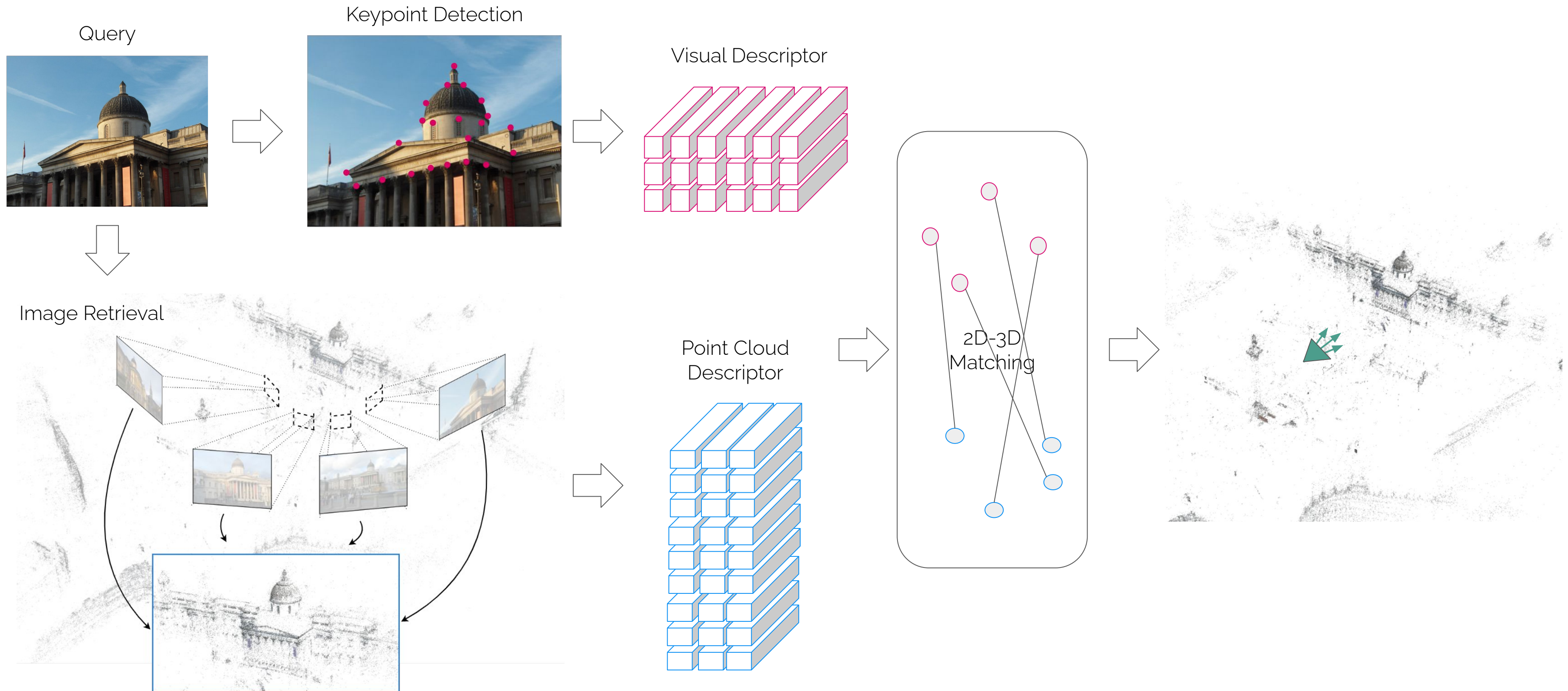Image Retrieval

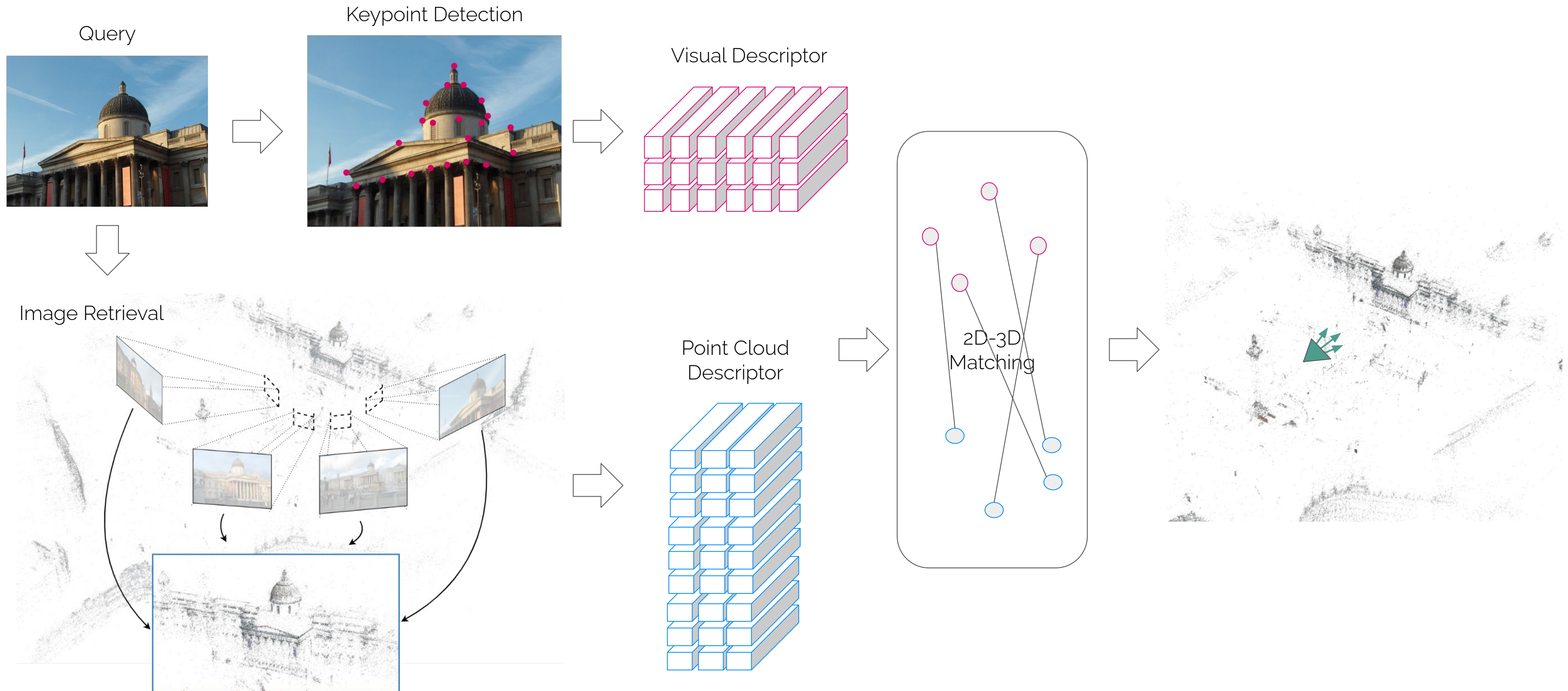# Classical Structure-based Localization

# Classical Structure-based Localization

# Classical Structure-based Localization



Query

Keypoint Detection

Visual Descriptor

Image Retrieval
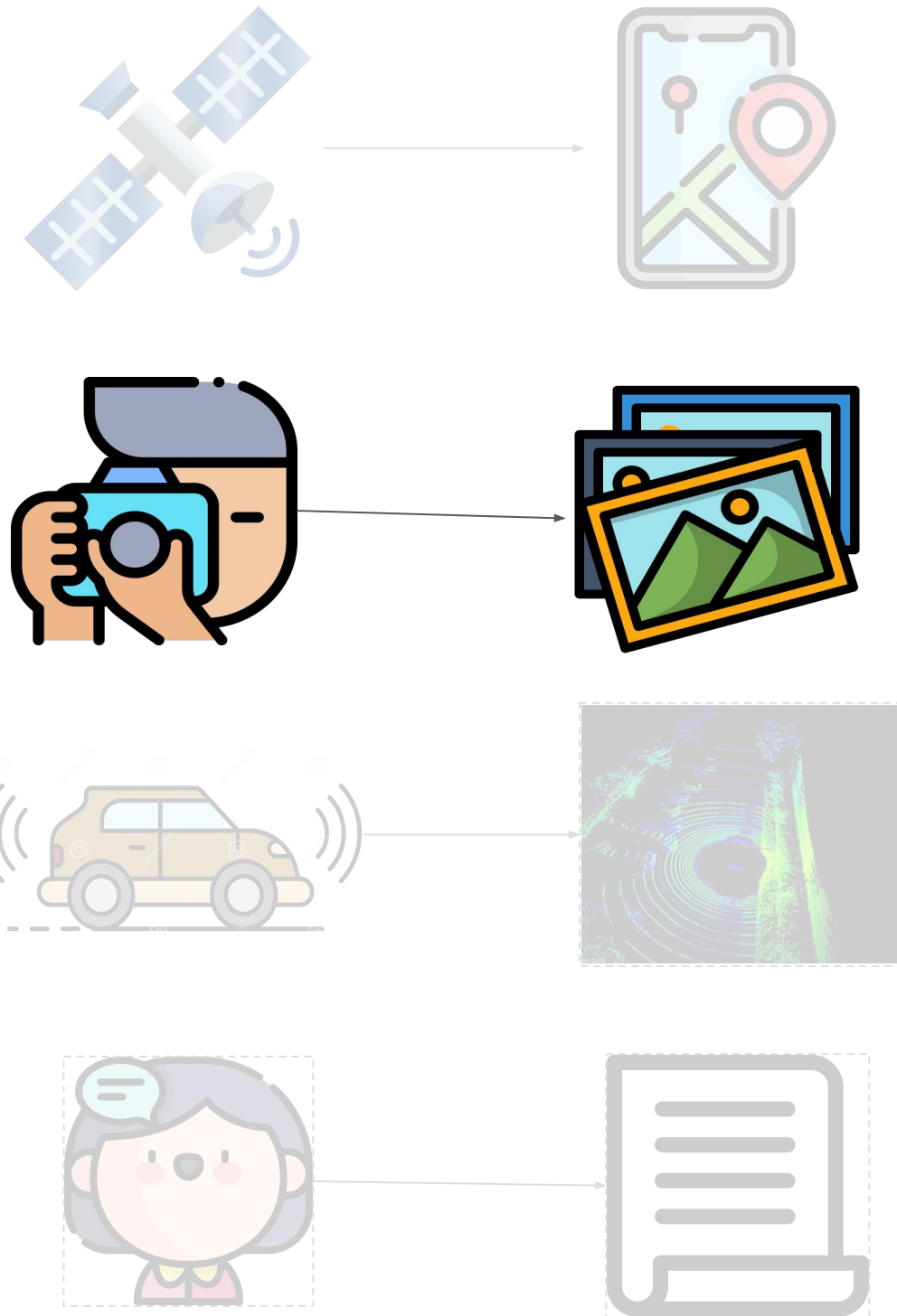
Point Cloud Descriptor

2D-3D Matching

# Classical Structure-based Localization

# Visual Localization



Query Data

Map Data

Localization System

Camera Pose

Orientation    Position
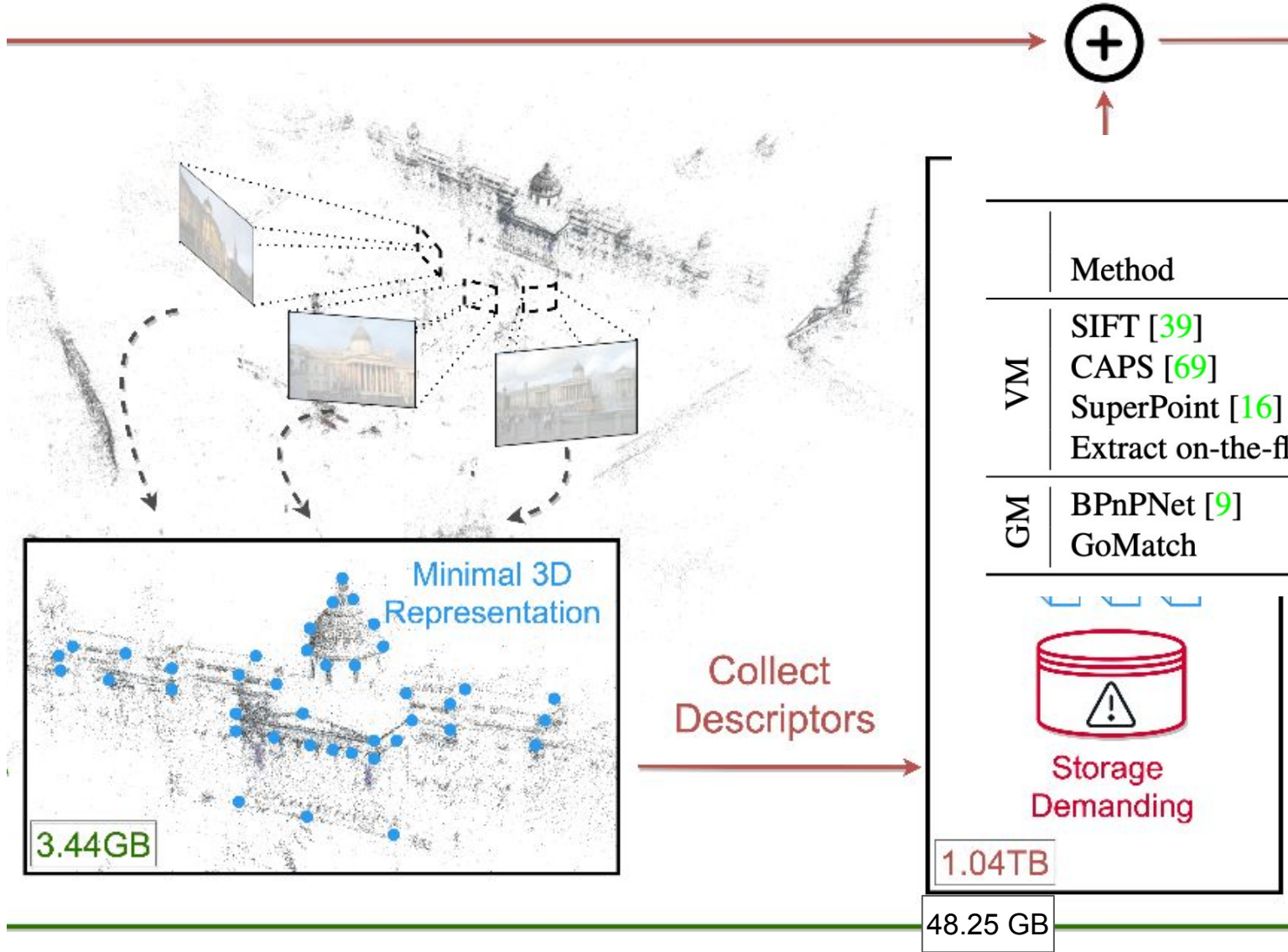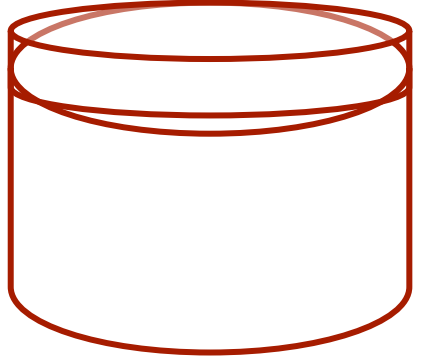
# AR/VR



https://blog.helpdocs.io/guidigo/



[Middelberg, Sattler, Untzelmann, Kobbelt, Scalable 6-DOF Localization on Mobile Devices, ECCV 2014]

# Storage Requirements



| | Method | Easy Maintenance | Privacy | Database Storage (GB, ↓) | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | Cameras (MB) | 3D | Raw Imgs | Descs | Total |
| VM | SIFT [39] | ✗ | ✗ | 15.73 | 3.44 | ✗ | 130.10 (uint8) | 133.33 |
| | CAPS [69] | ✗ | ✗ | 15.73 | 3.44 | ✗ | 520.38 (fp32) | 523.83 |
| | SuperPoint [16] | ✗ | ✗ | 15.73 | 3.44 | ✗ | 1040.76 (fp32) | 1044.21 |
| | Extract on-the-fly | ✗ | ✗ | 15.73 | 3.44 | 157.84 | ✗ | 161.29 |
| GM | BPnPNet [9] | ✓ | ✓ | 15.73 | 3.44 | ✗ | ✗ | **3.45** |
| | GoMatch | ✓ | ✓ | | | | | |

Minimal 3D Representation

Collect Descriptors

Storage Demanding

3.44GB

1.04TB

48.25 GB

# Privacy Challenge



(a) SfM point cloud (top view)    (b) Projected 3D points    (c) Synthesized Image    (d) Original Image

nD Input Tensor

z RGB    SIFT descriptor

encoder    decoder    conv. layers

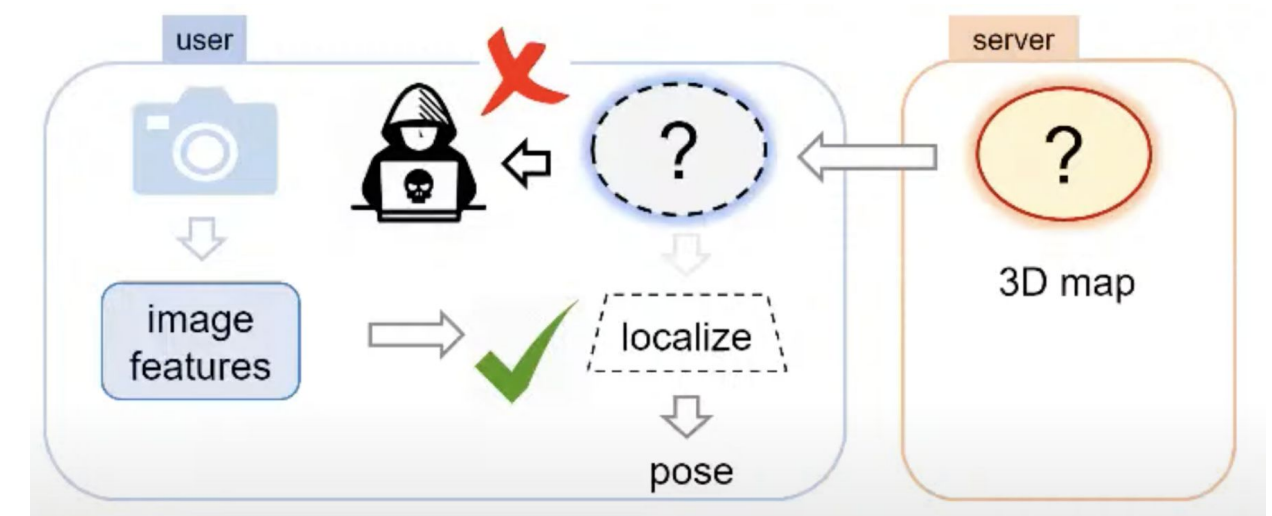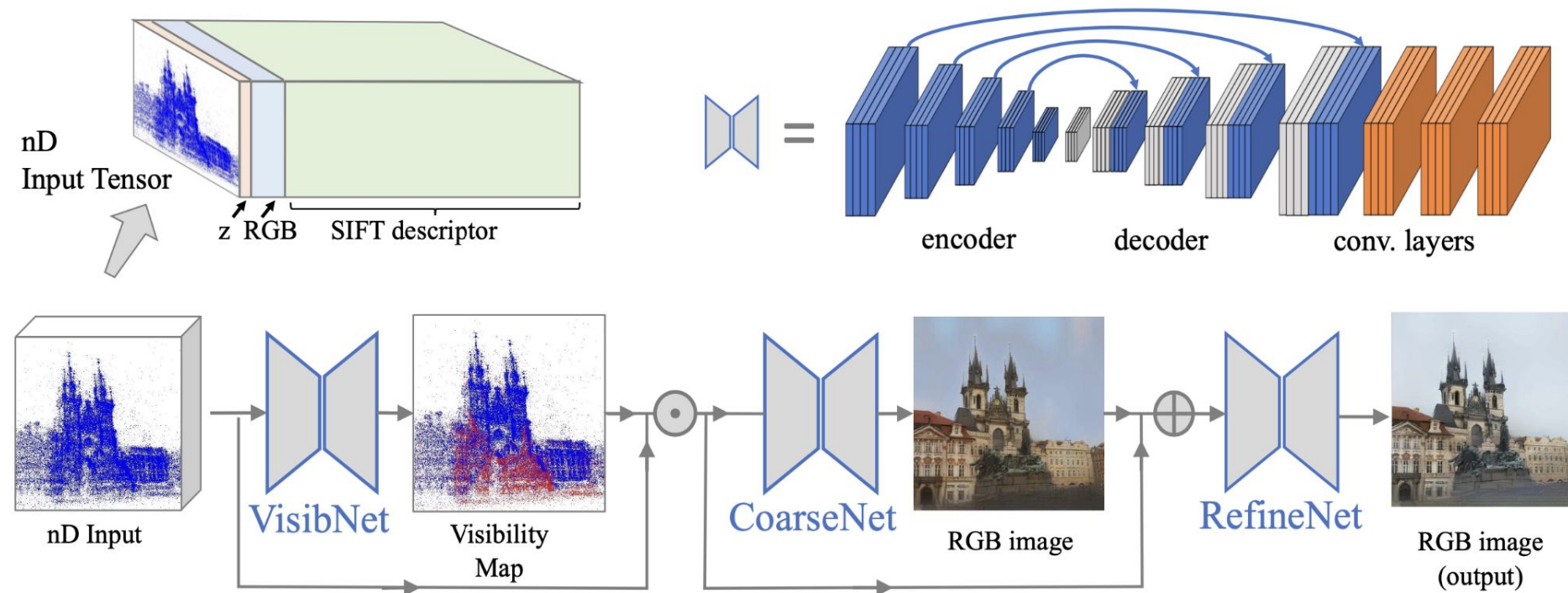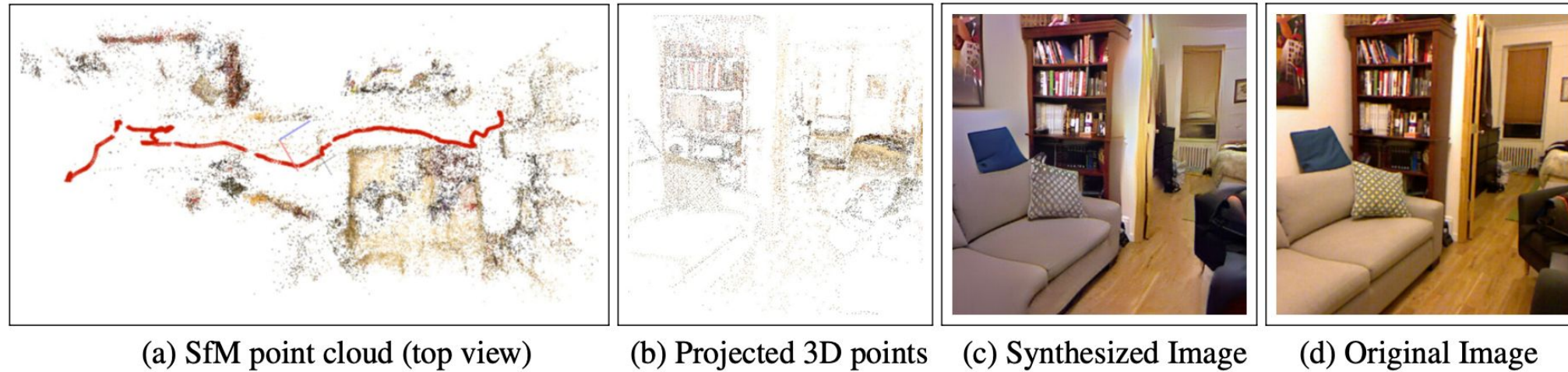nD Input    VisibNet    Visibility Map    CoarseNet    RGB image    RefineNet    RGB image (output)

Francesco, Pittaluga, et al Revealing Scenes by Inverting Structure From Motion Reconstructions. CVPR19



**Man-in-the-middle Attack**

# Privacy Challenge



(a) SfM point cloud (top view)  (b) Projected 3D points  (c) Synthesized Image  (d) Original Image

nD Input Tensor

z RGB  SIFT descriptor

encoder  decoder  conv. layers

nD Input  VisibNet  Visibility Map  CoarseNet  RGB image  RefineNet  RGB image (output)
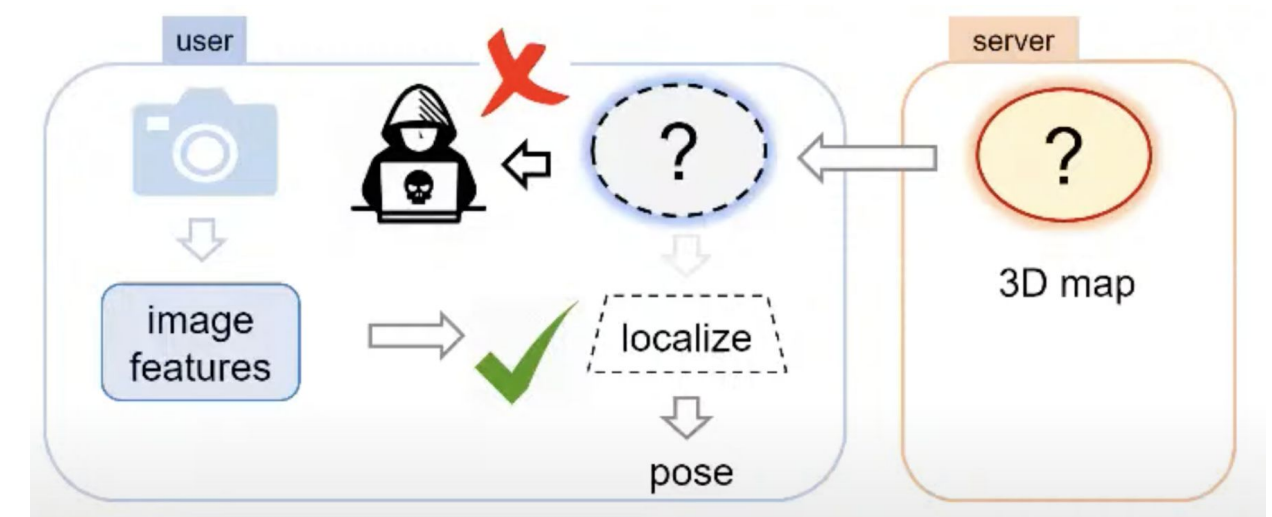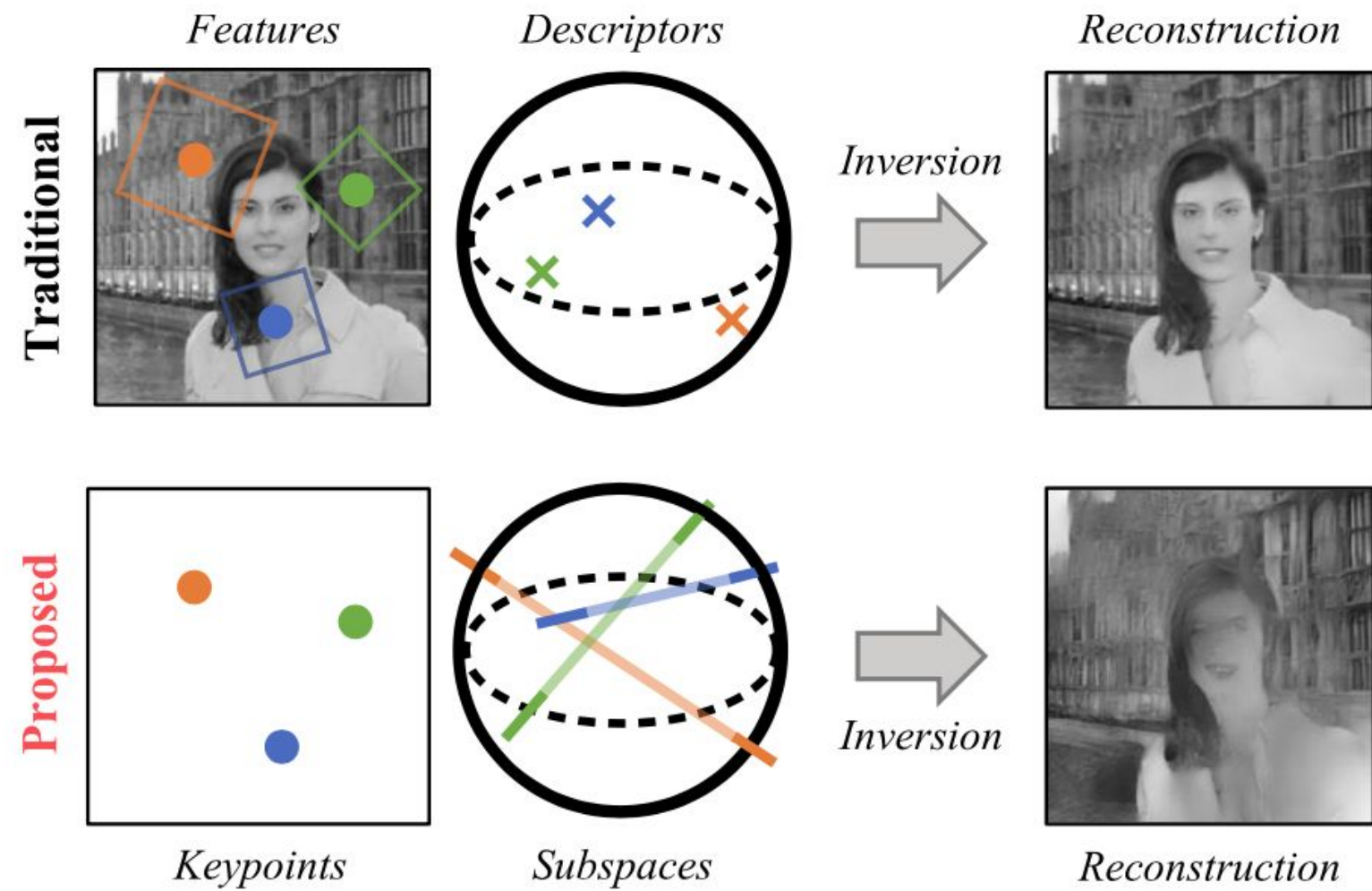
Francesco Pittaluga, Sanjeev J.Koppal, Sing Bing Kang, and Sudipta N Sinha. Revealing Scenes by Inverting Structure From Motion Reconstructions. CVPR19

user  cloud  map

features  localization

pose

**Man-in-the-middle Attack**

user  server

image features  localize  3D map
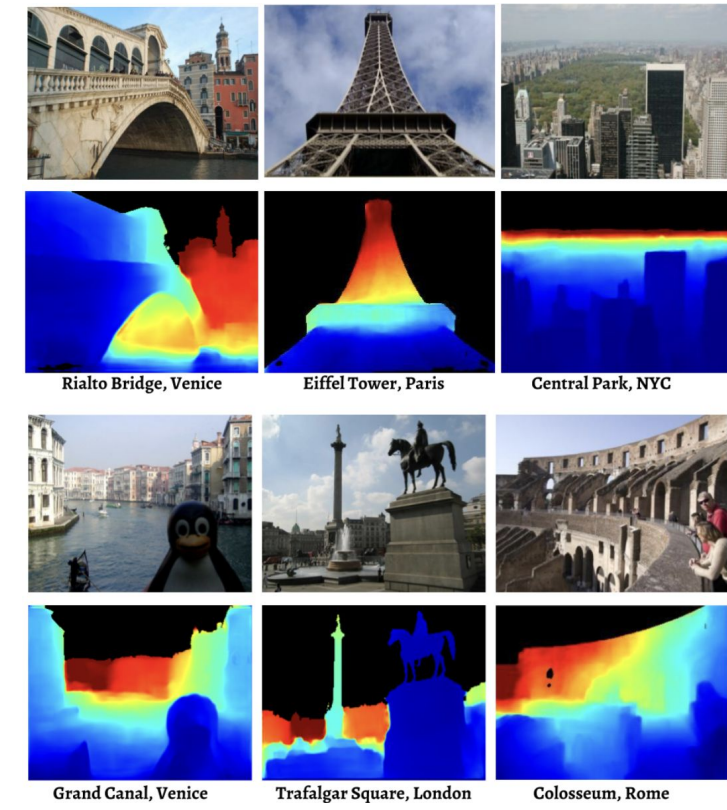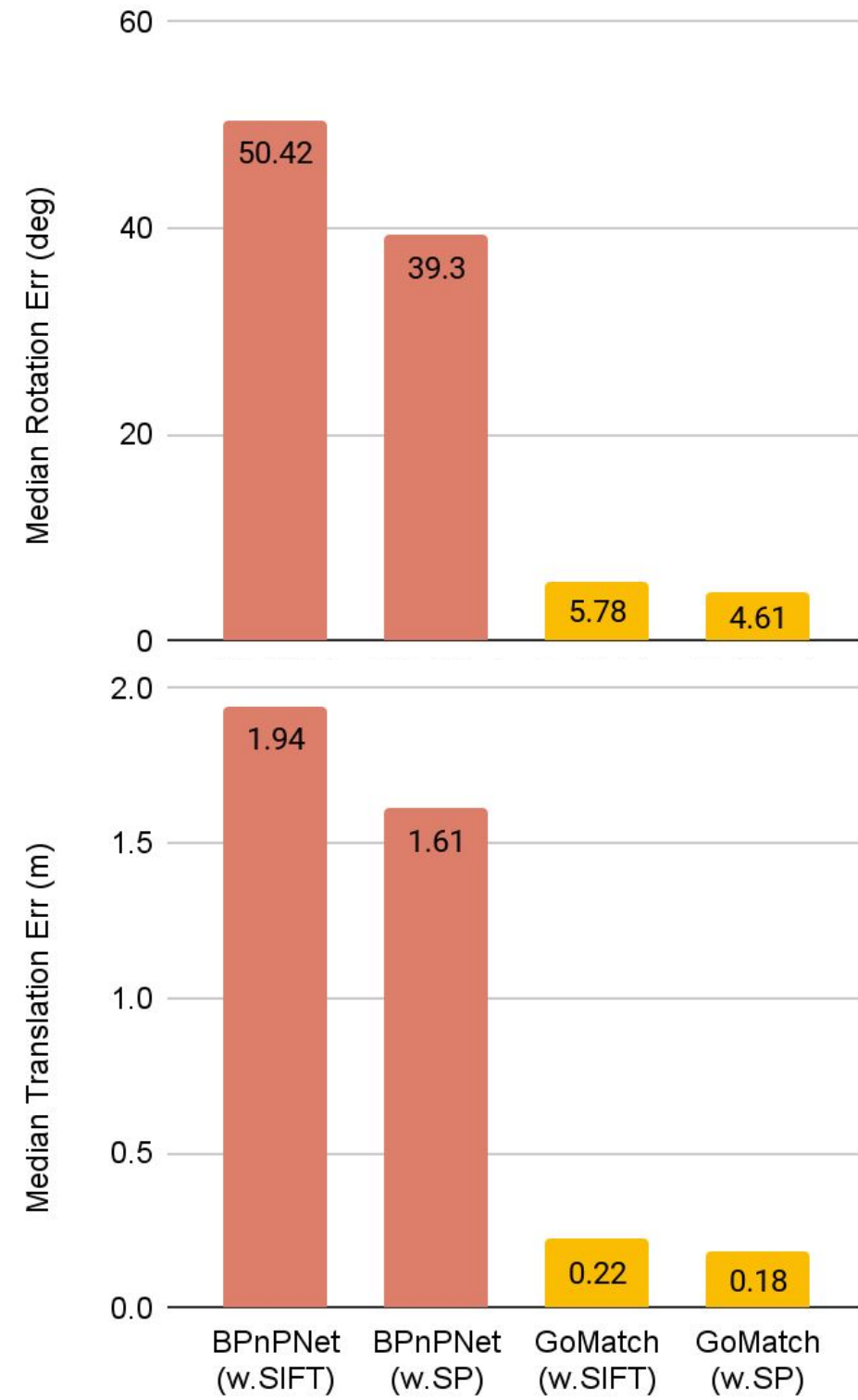
pose

# Privacy Challenge



Dusmanu, Mihai, et al. "Privacy-preserving image features via adversarial affine subspace embeddings." CVPR21.



Ng, Tony, et al. "NinjaDesc: Content-Concealing Visual Descriptors via Adversarial Learning." CVPR22

# Generalization



7Scenes



Rialto Bridge, Venice — Eiffel Tower, Paris — Central Park, NYC

Grand Canal, Venice — Trafalgar Square, London — Colosseum, Rome

# Evaluation