

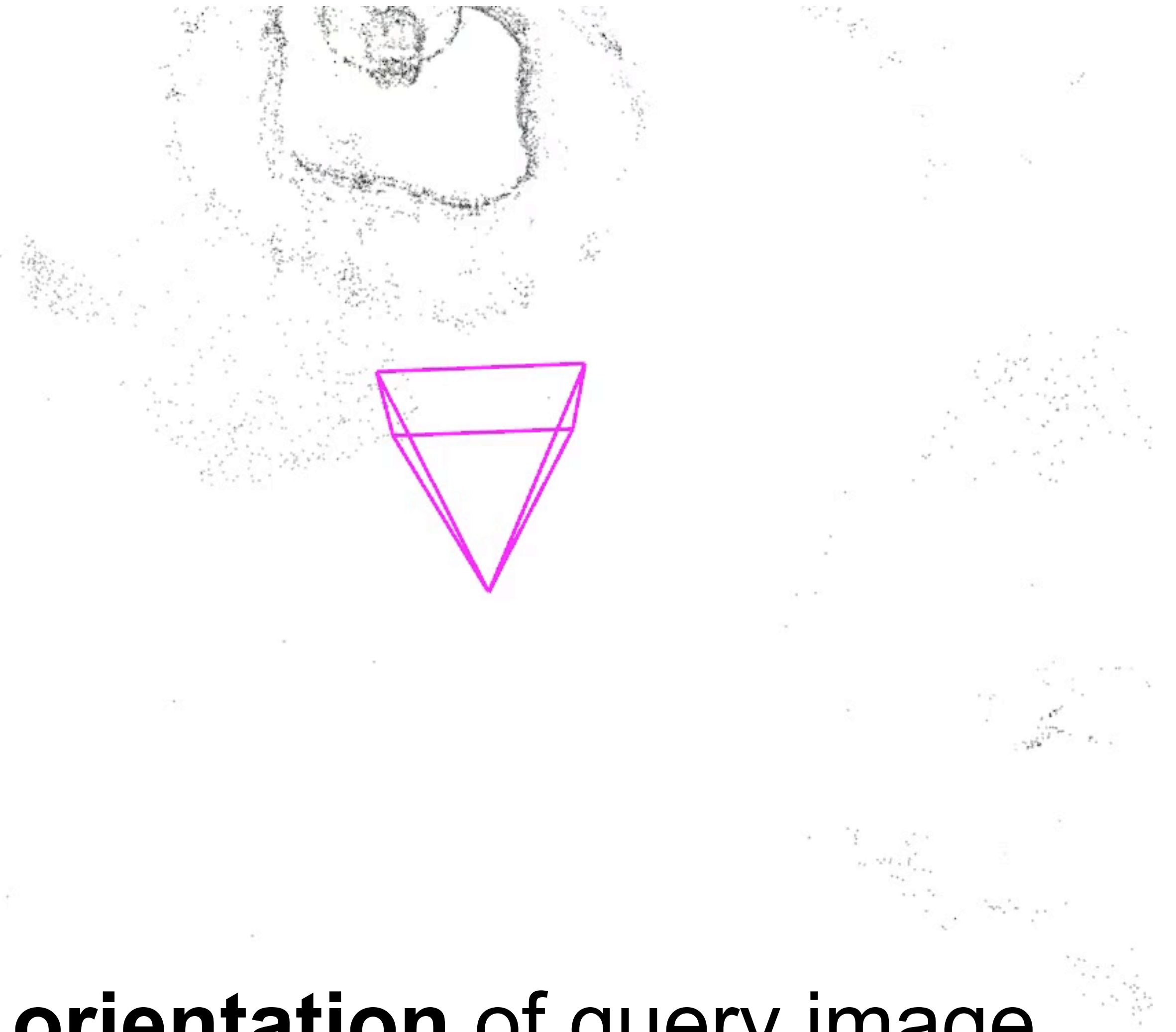
Scene Representations for Visual Localization

Torsten Sattler

Czech Institute of Informatics, Robotics and Cybernetics
Czech Technical University in Prague

Torsten Sattler

The Visual Localization Problem

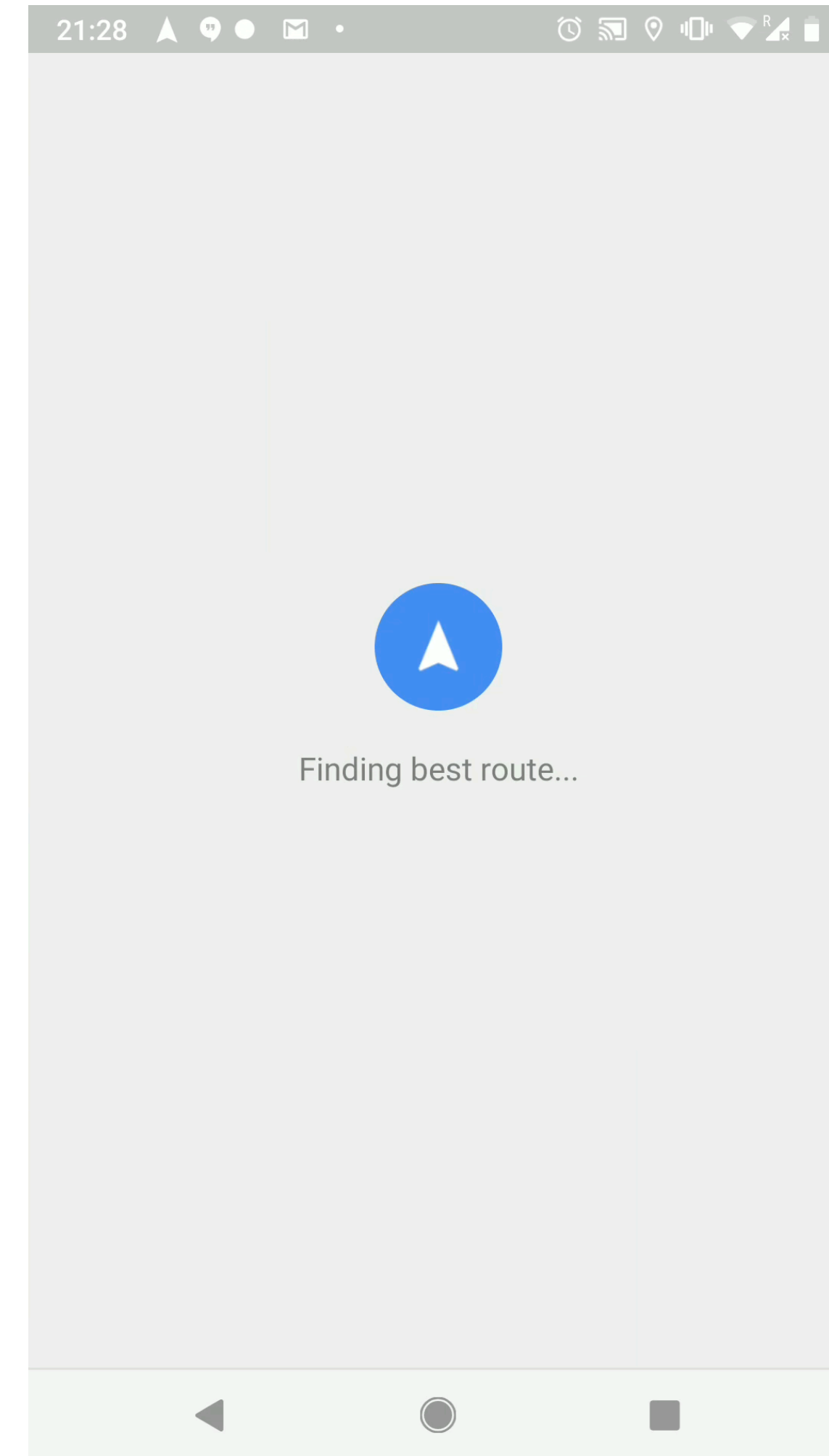


Compute **exact position and orientation** of query image

Applications: Augmented Reality

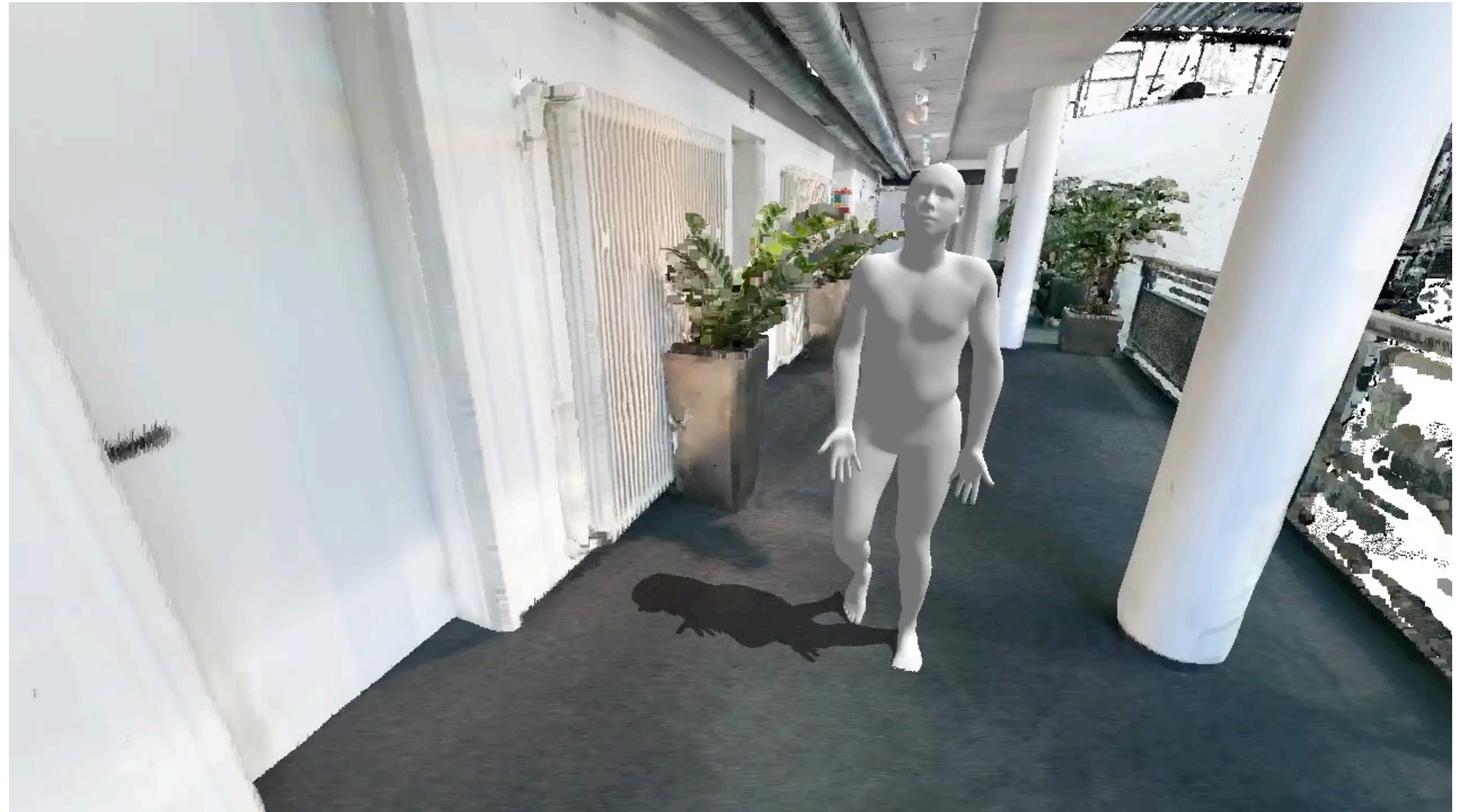


[Middelberg, Sattler, Untzelmann, Kobbelt, Scalable 6-DOF Localization on Mobile Devices, ECCV 2014]



AR navigation in Google Maps

Applications: Performance Capture



slide credit: Vladimir Guzov, Aymen Mir

[Guzov*, Mir*, Sattler, Pons-Moll, Human POSEitioning System (HPS): 3D Human Pose Estimation and Self-localization in Large Scenes from Body-Mounted Sensors, CVPR 2021]

Torsten Sattler

Applications: Visual Localization for Modeling Interactions



slide credit: Vladimir Guzov

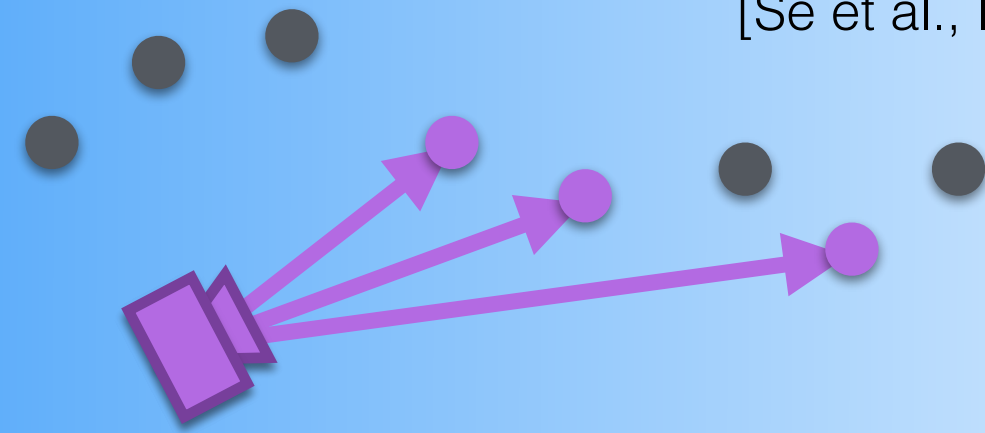
[Guzov, Chibane, Marin, He, Sattler, Pons-Moll, Interaction Replica: Tracking human–object interaction and scene changes from human motion, arXiv 2023]

Visual Localization - A Taxonomy

3D structure-based representation

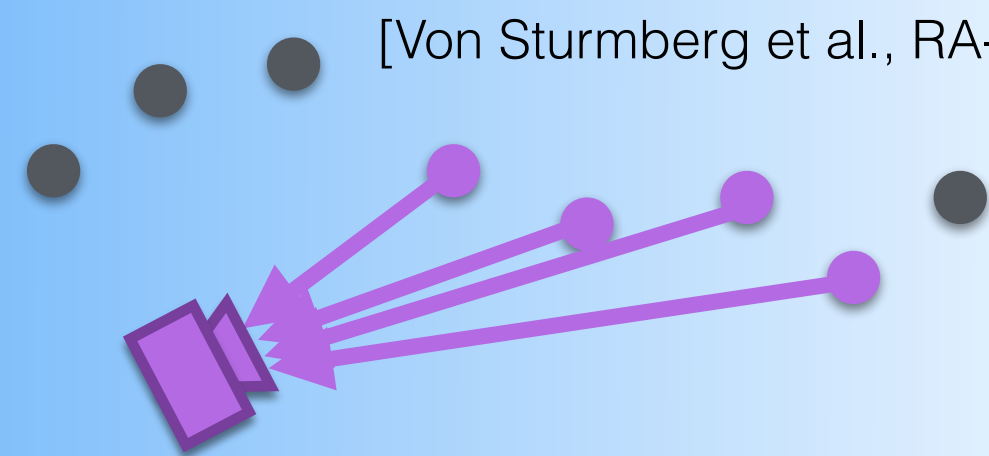
feature-based localization

[Se et al., IROS'02]



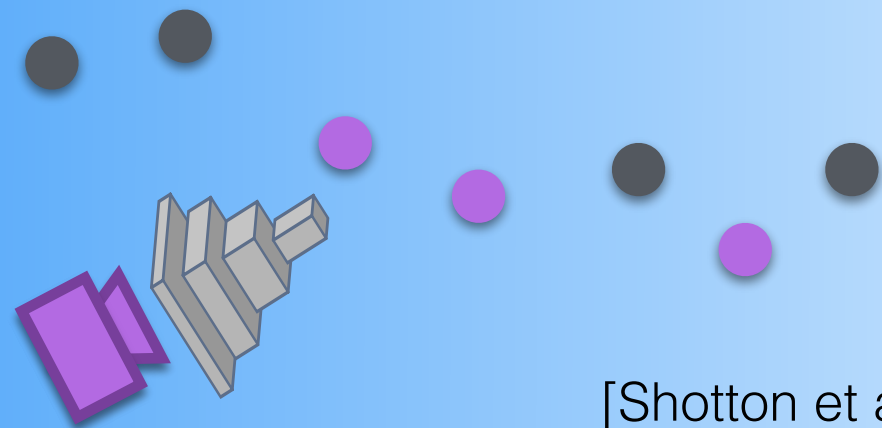
pose refinement

[Von Sturmberg et al., RA-L'20]



scene coordinate regression

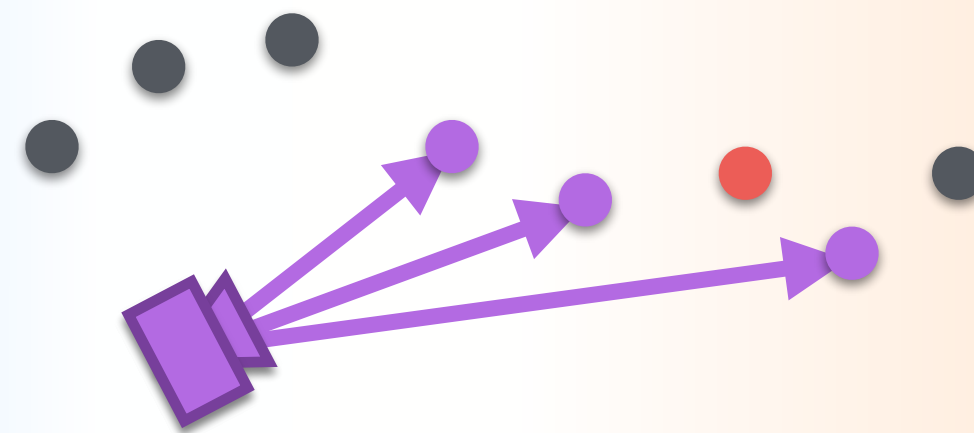
[Shotton et al., CVPR'13]



hybrid

hierarchical localization

[Irschara et al., CVPR'09]



hybrid pose estimation

[Camposeco et al., ICCV'21]

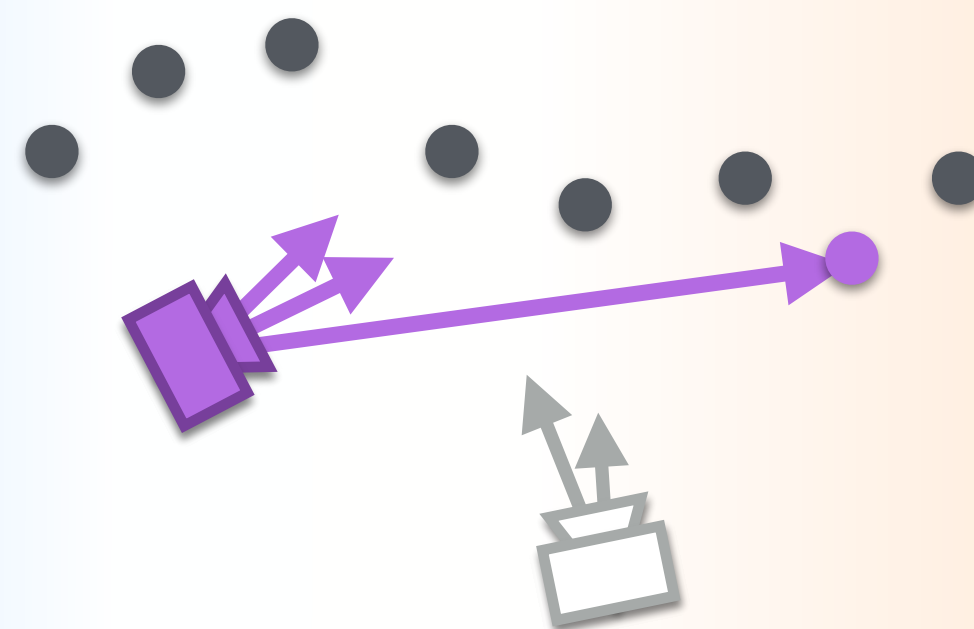


image-based representation

pose triangulation

[Zhang & Kosecka, 3DPVT'06]

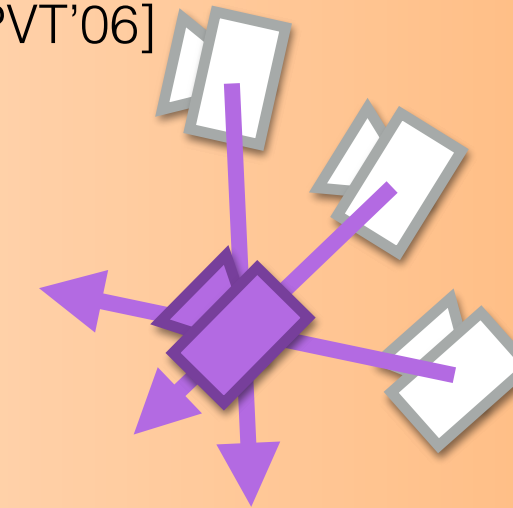
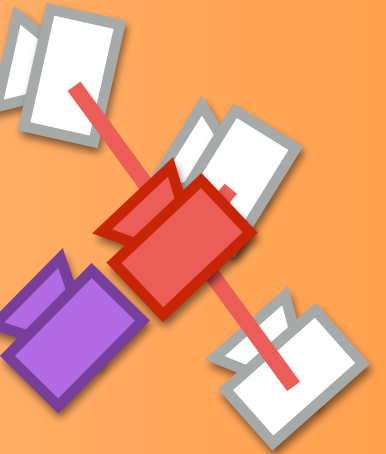


image retrieval



pose interpolation

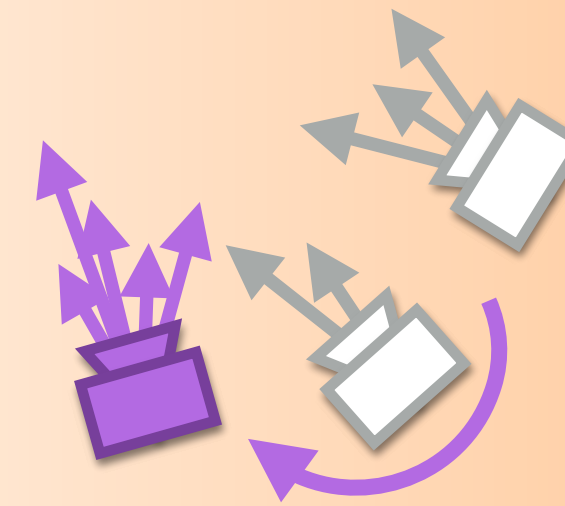
[Torii et al., ICCV'11]



semi-generalized relative pose / homography

[Zheng & Wu, ICCV'15]

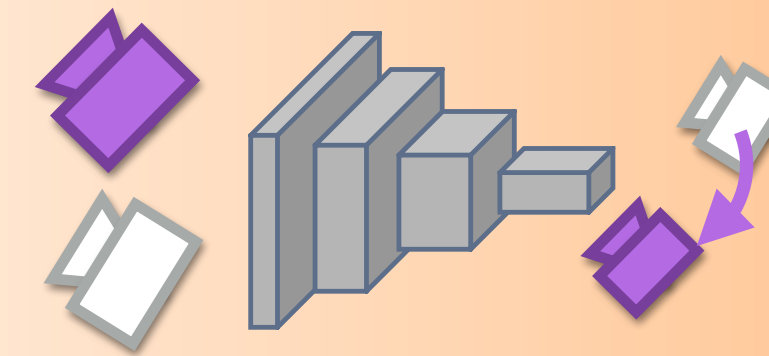
[Bhayani et al., ICCV'21]



relative pose regression

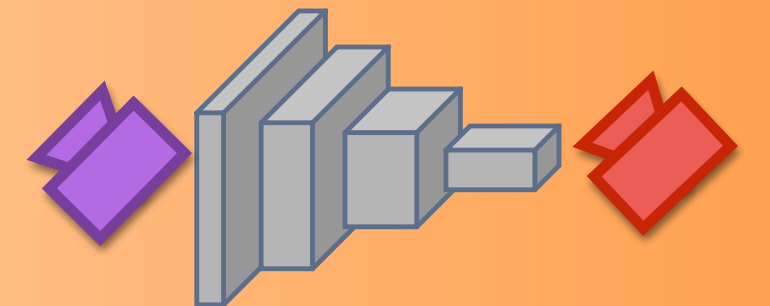
[Laskar et al., CVPRW'17]

[Balntas et al., ECCV'18]

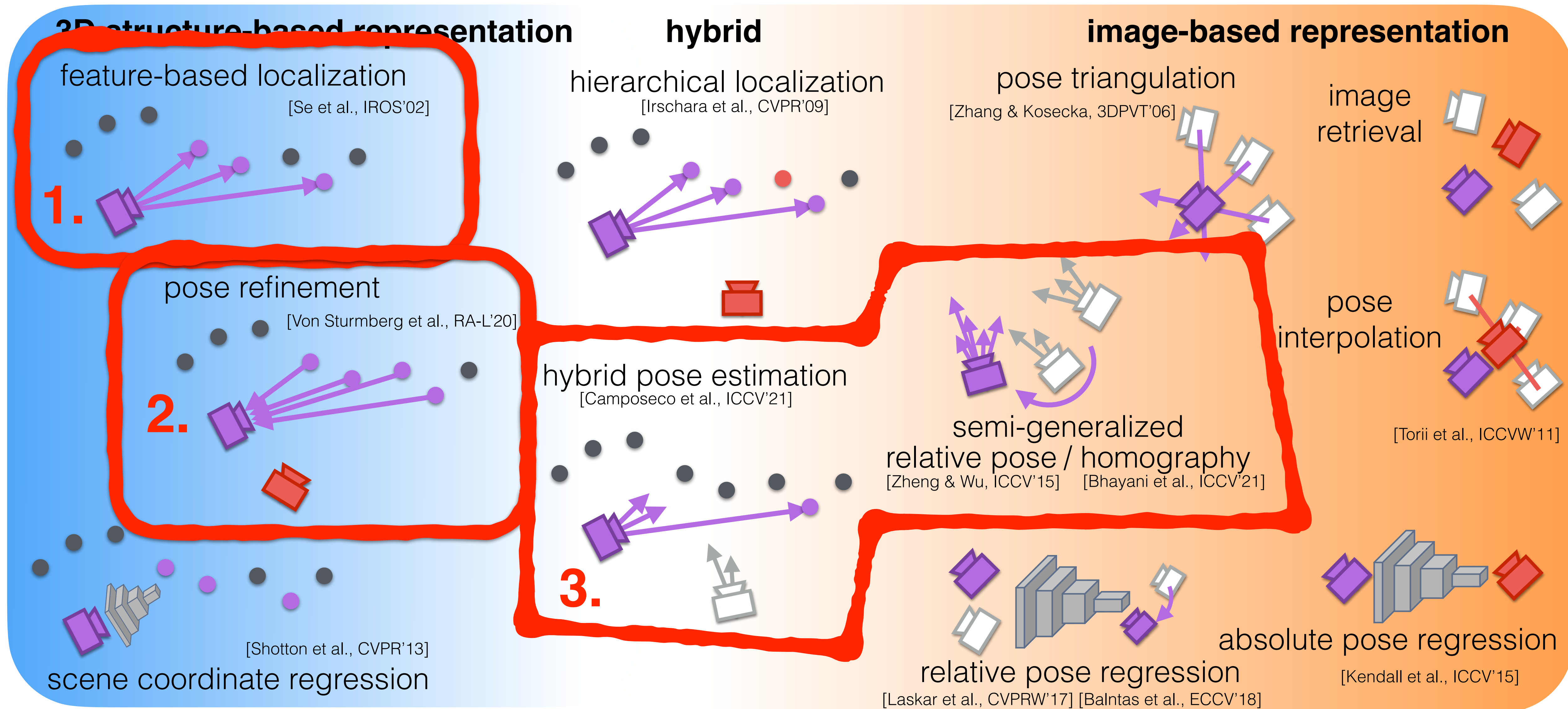


absolute pose regression

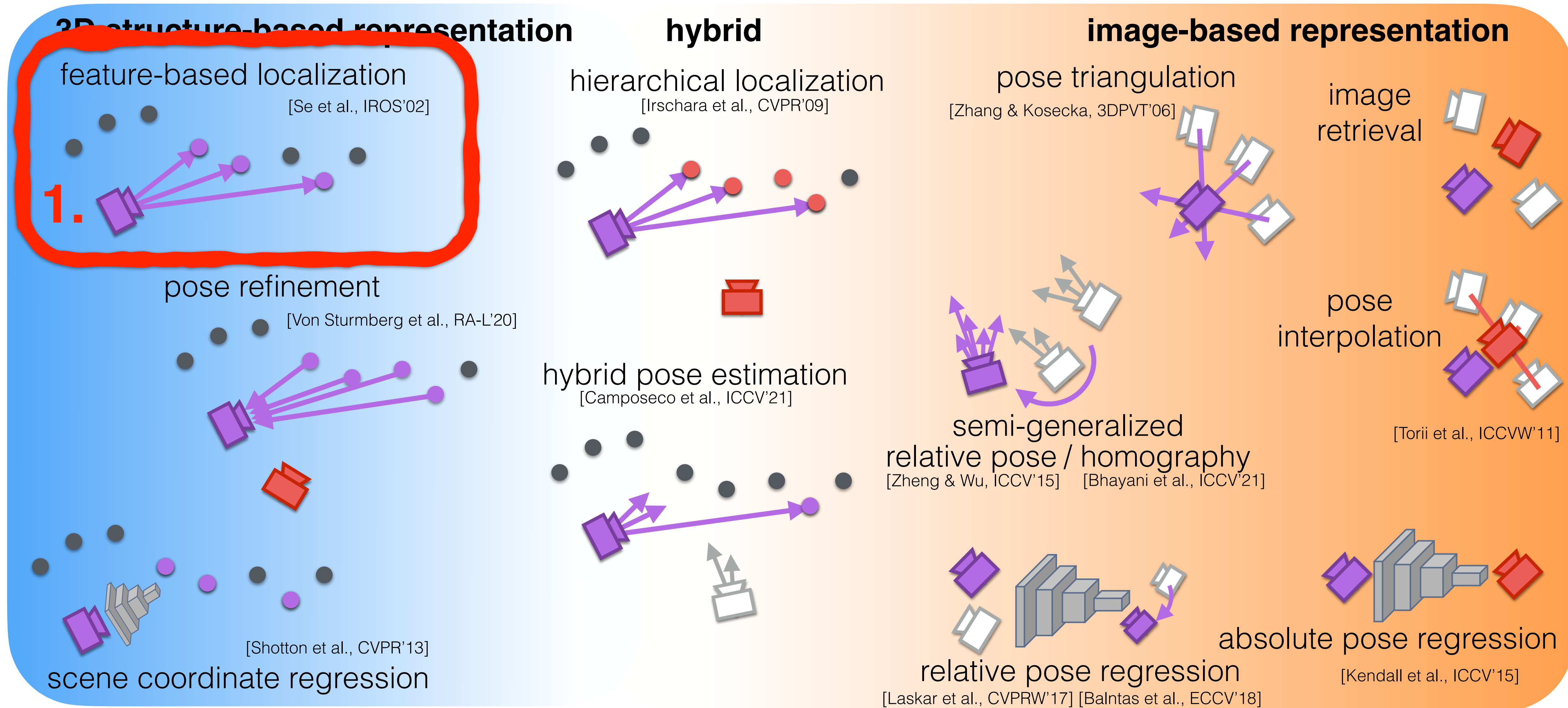
[Kendall et al., ICCV'15]



Visual Localization - A Taxonomy



Visual Localization - A Taxonomy



Classical Representation: SfM Point Clouds

3D point triangulated from ≥ 2 images:
3D position + local feature descriptors



For new query image:
Establish 2D-3D matches via
feature matching



Pose estimation from 2D-3D
correspondences



Classical Representation: SfM and Clouds

SIFT
(1999)

3D point cloud regulated from ≥ 2 images:
3D position + local feature descriptors

kd-trees (1975),
BoW (2003/4)

For new image:
Establish 2D-3D matches via
feature matching

Pose estimation from 2D-3D
correspondences

P3P (≤ 1773)
RANSAC (1981)

Classical Representation: SfM Point Clouds



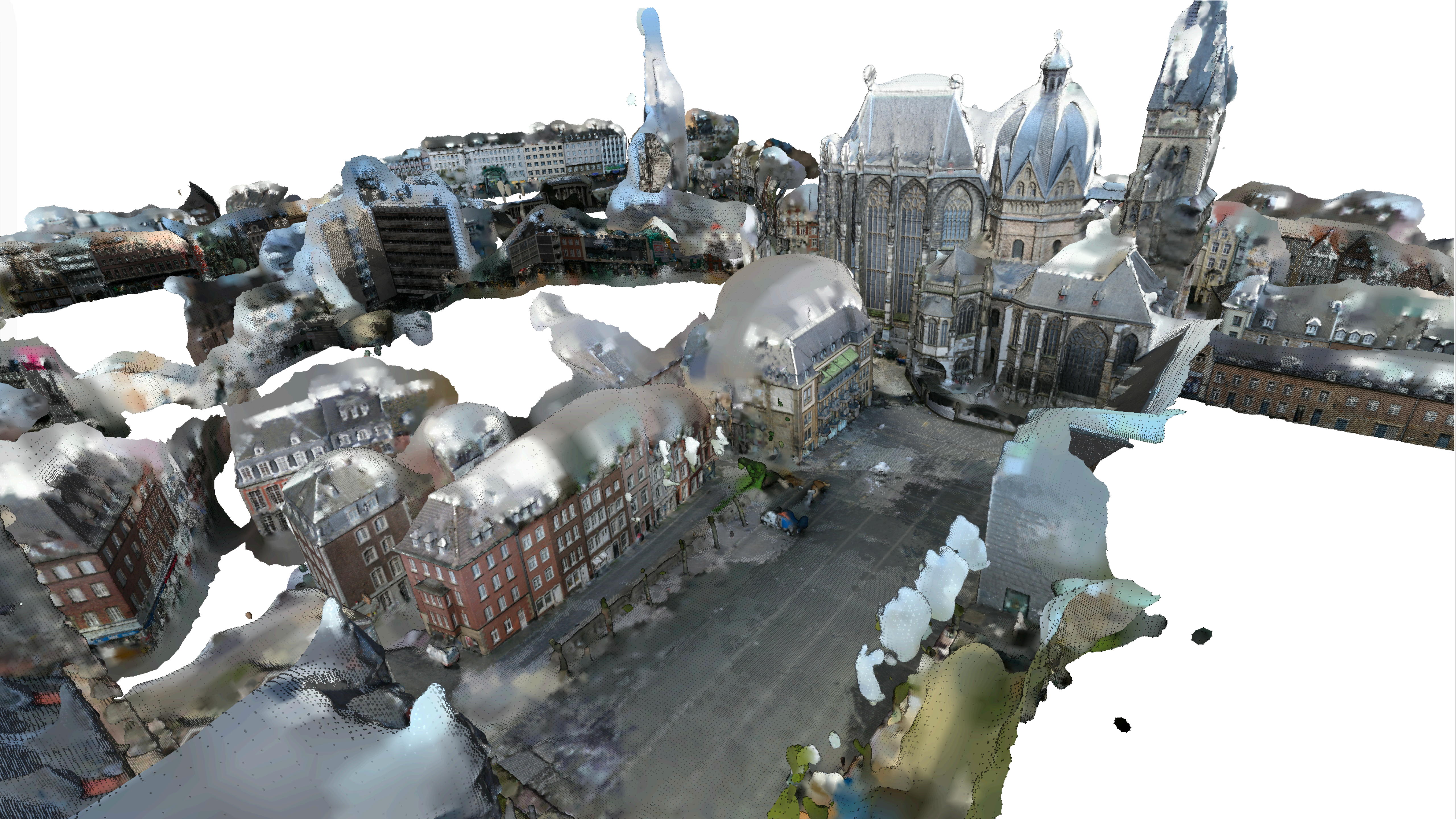
Classical Representation: SfM Point Clouds

Advantages:

- ✓ Efficient
- ✓ Scalable
- ✓ Quite robust to condition changes
- ✓ Easily compressible

Disadvantages:

- ✗ Specialized & sparse representation
- ✗ Needs to be recomputed when changing features



SfM-based vs. Mesh-based Localization

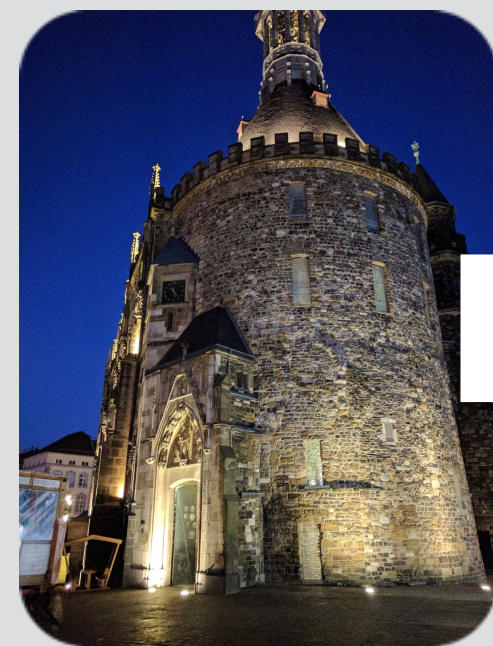
offline (input: posed images)

Extract features from database images → Match features between database images → Triangulate scene structure

offline (input: posed images)

3D reconstruction

online



Find N most similar database images

Feature matching between query and top-N database images

Pose estimation from 2D-3D matches (P3P-LO-RANSAC)

2D-3D matches from associations between database features and 3D points

online



Find N most similar database images

Feature matching between query and top-N database images

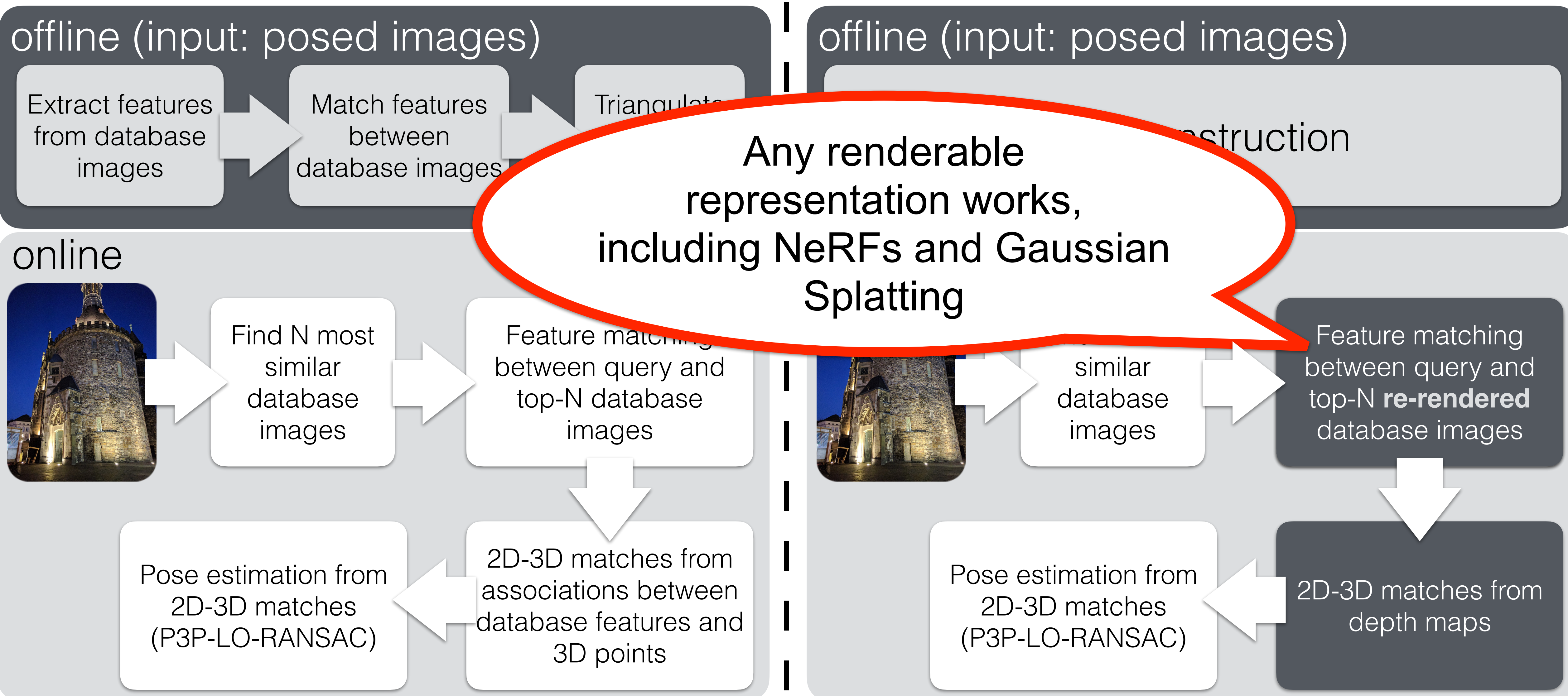
Pose estimation from 2D-3D matches (P3P-LO-RANSAC)

2D-3D matches from depth maps

[Panek, Kukelova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

SfM-based vs. Mesh-based Localization



[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

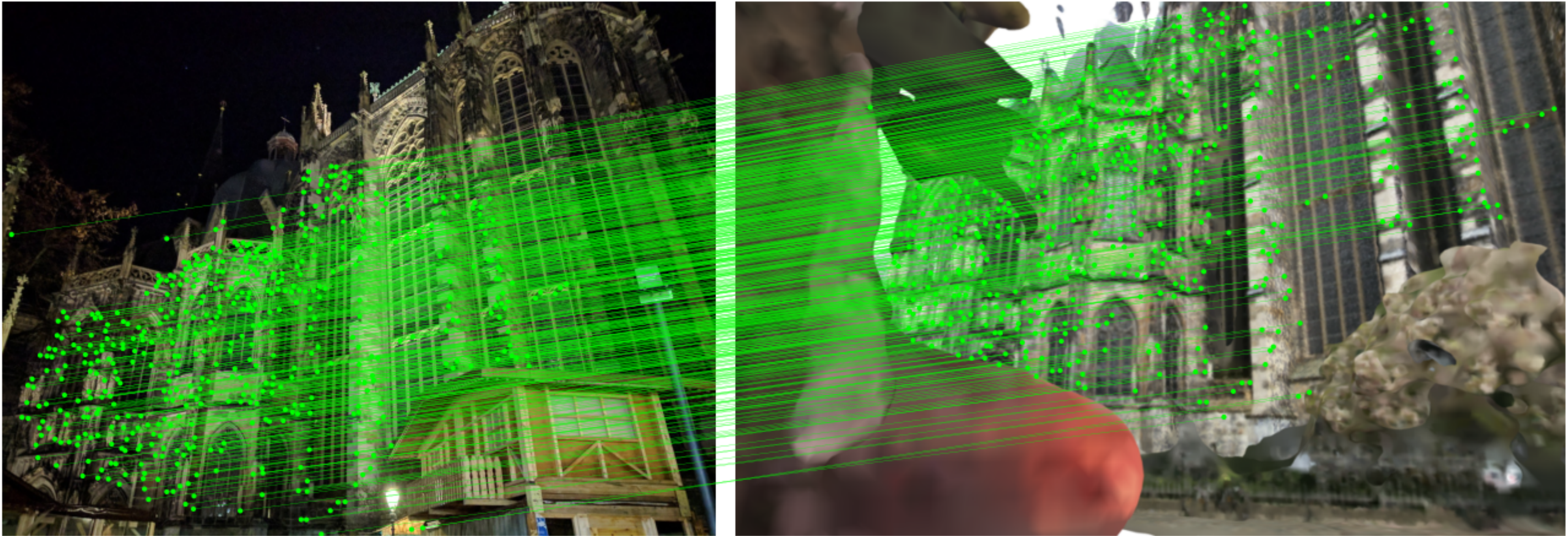
Matching Against Rendered Images



[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

Matching Against Rendered Images

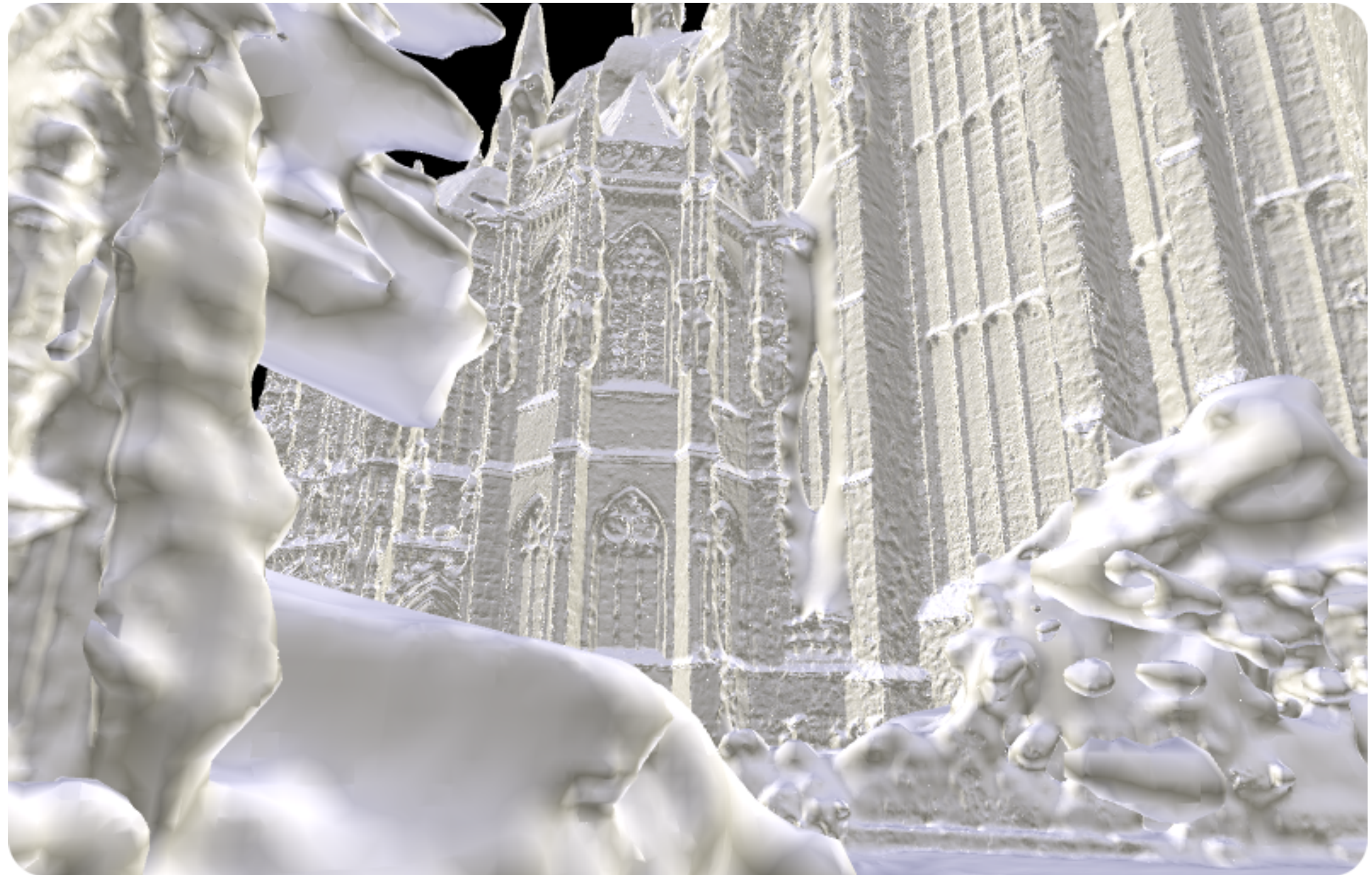


ALIKED features [Zhao et al., IEEE TIM 2023] with LightGlue [Lindemberger et al., ICCV 2023]
matcher not trained on renderings

[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

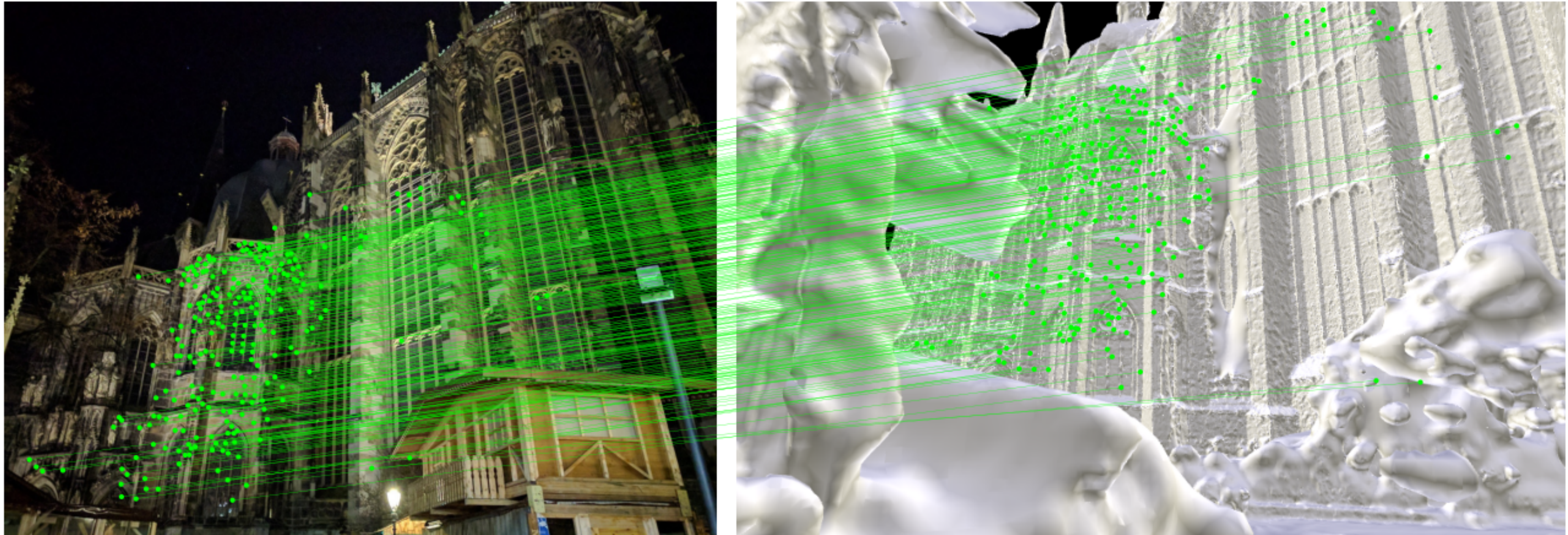
Matching Against Raw Geometry



[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

Matching Against Raw Geometry

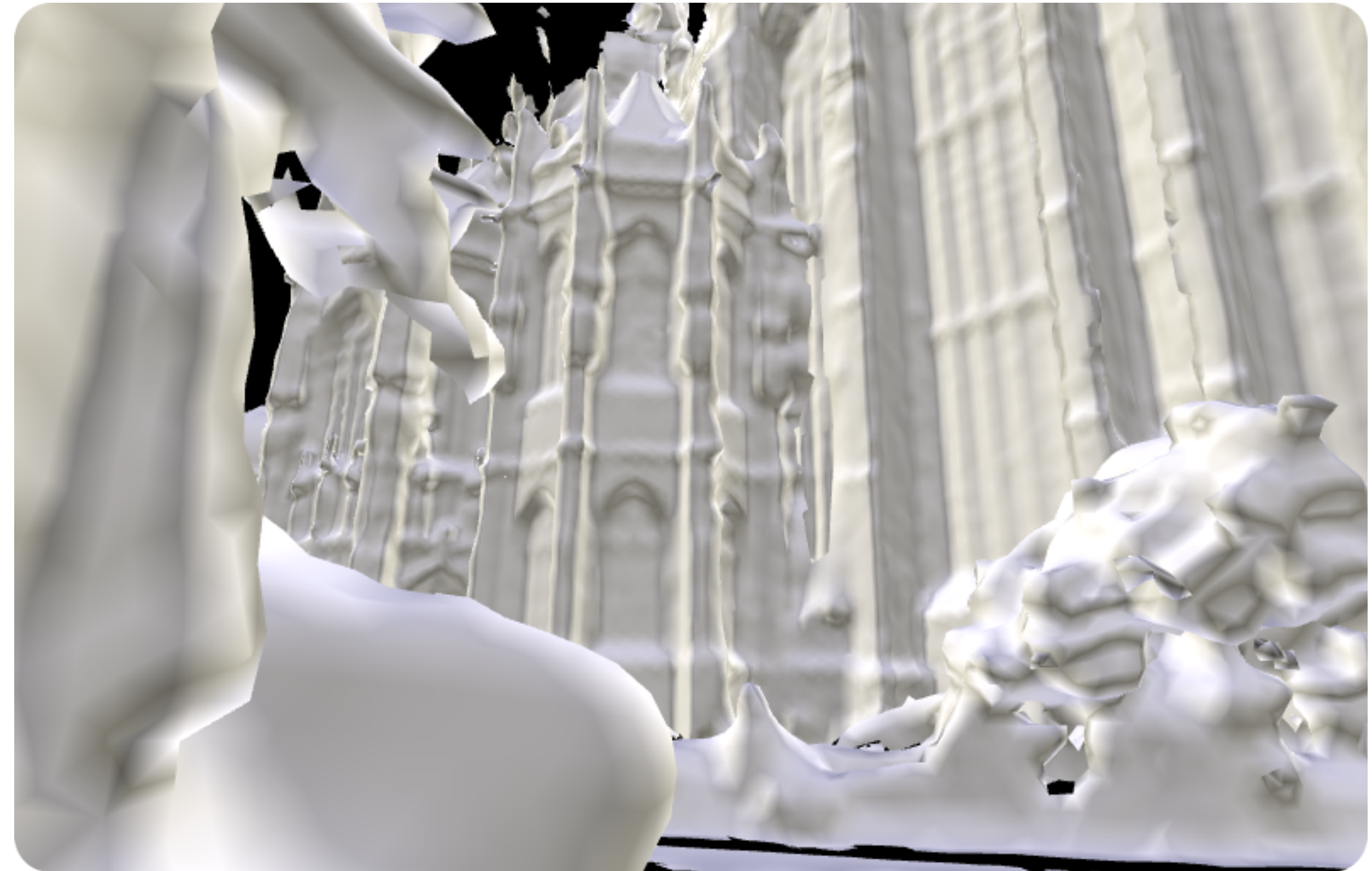


ALIKED features [Zhao et al., IEEE TIM 2023] with LightGlue [Lindemberger et al., ICCV 2023]
matcher not trained on renderings

[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

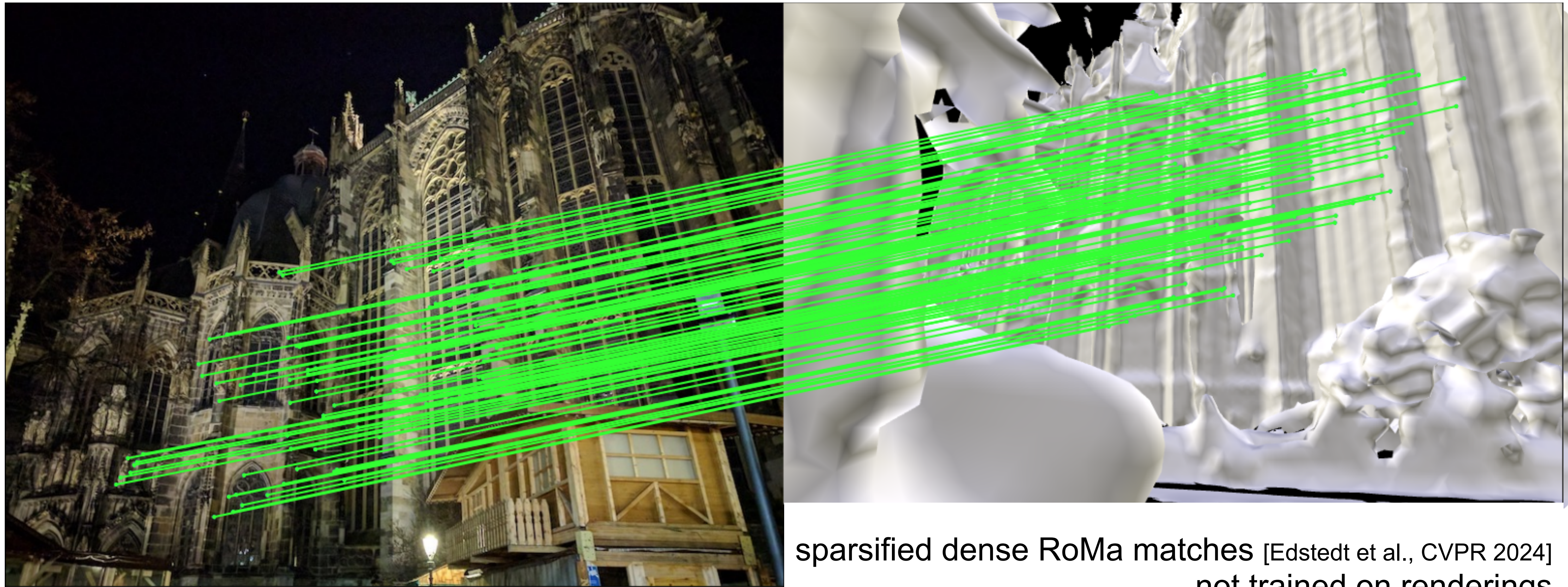
Matching Against Raw Geometry



[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

Torsten Sattler

Matching Against Raw Geometry



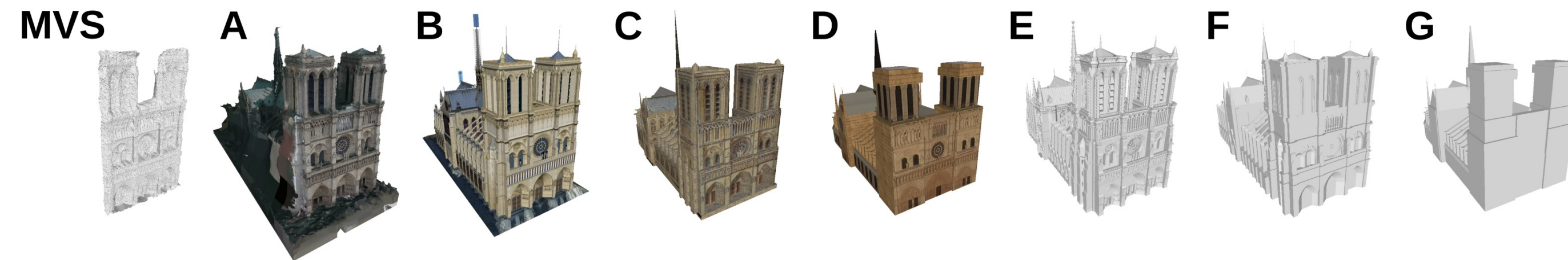
sparsified dense RoMa matches [Edstedt et al., CVPR 2024]
not trained on renderings

[Panek, Kukulova, Sattler, MeshLoc: Mesh-Based Visual Localization, ECCV 2022]

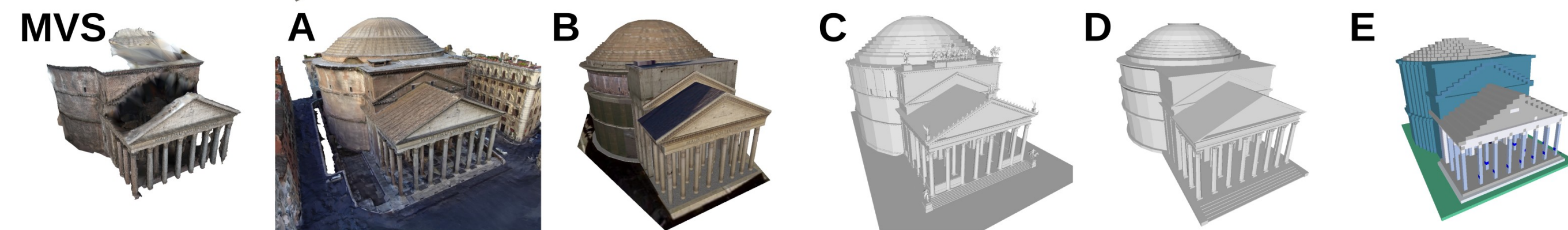
Torsten Sattler

Visual Localization using Internet Models

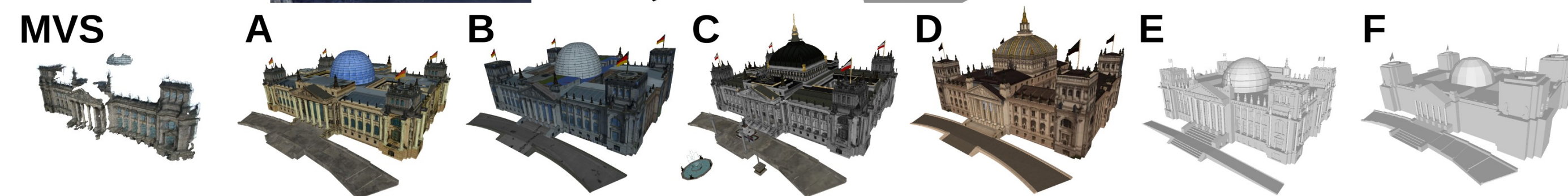
Notre Dame



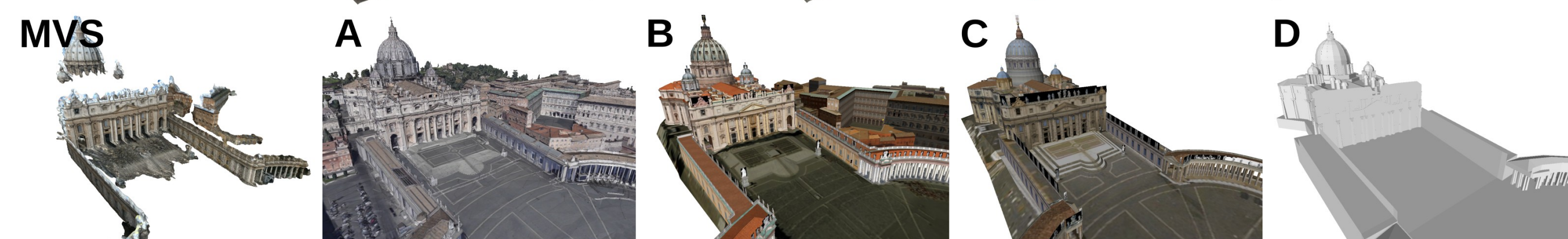
Pantheon



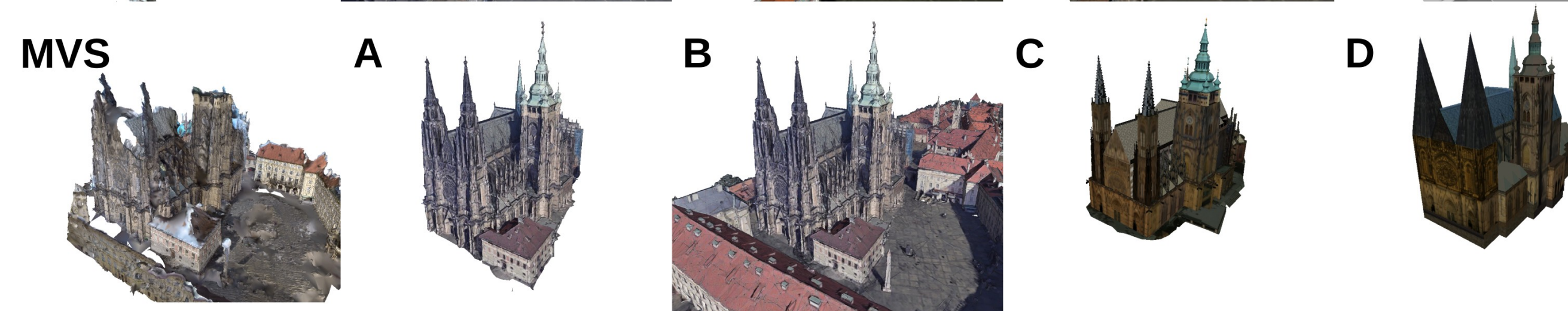
Reichstag



St. Peter's Square



St. Vitus Cathedral



[Panek, Kukelova, Sattler, Visual Localization using Imperfect 3D Models from the Internet, CVPR 2023]

Torsten Sattler

Benchmarking Visual Localization using Internet Models

VISUAL LOCALIZATION ON 3D MESH MODELS

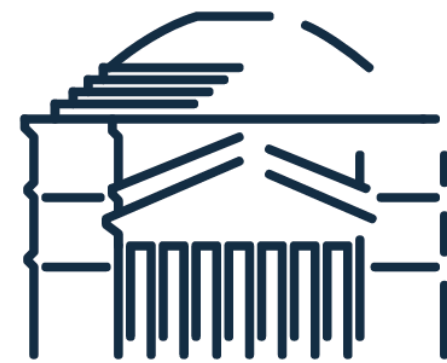
BENCHMARK

DATASETS

INFO



Notre Dame
(Paris, France)



Pantheon
(Rome, Italy)



Reichstag
(Berlin, Germany)



St. Peter's Square
(Vatican)



St. Vitus Cathedral
(Prague, Czech Republic)



v-pnk.github.io/cadloc

slide credit: Vojtech Panek

[Panek, Kukelova, Sattler, Visual Localization using Imperfect 3D Models from the Internet, CVPR 2023]

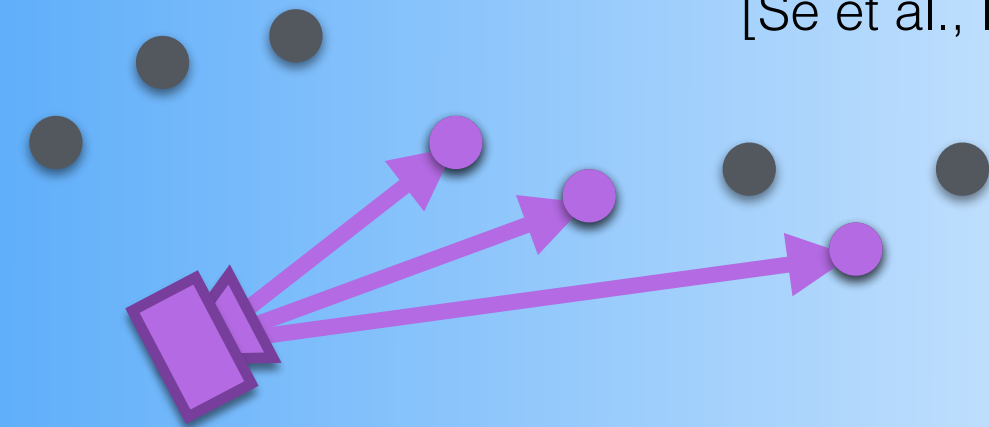
Torsten Sattler

Visual Localization - A Taxonomy

3D structure-based representation

feature-based localization

[Se et al., IROS'02]



pose refinement

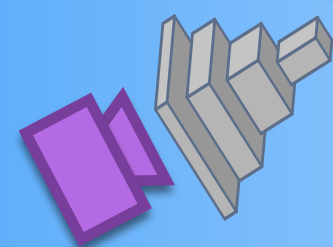
[Von Sturmberg et al., RA-L'20]

2.



scene coordinate regression

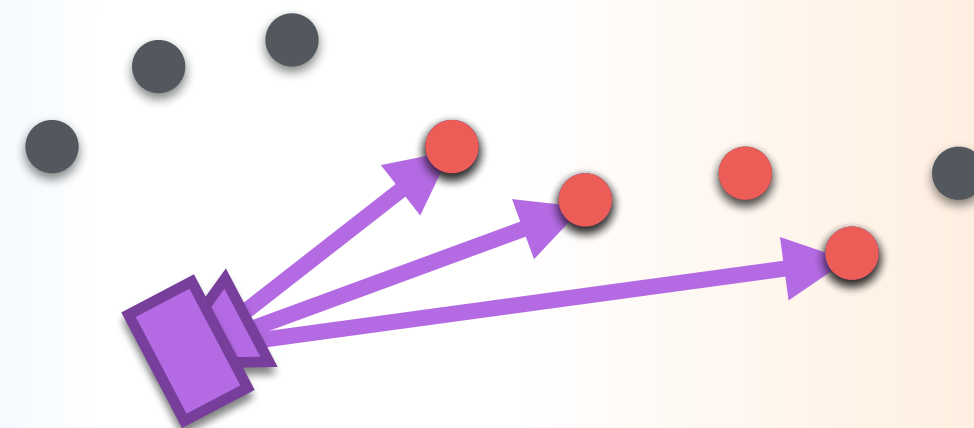
[Shotton et al., CVPR'13]



hybrid

hierarchical localization

[Irschara et al., CVPR'09]



hybrid pose estimation

[Camposco et al., ICCV'21]

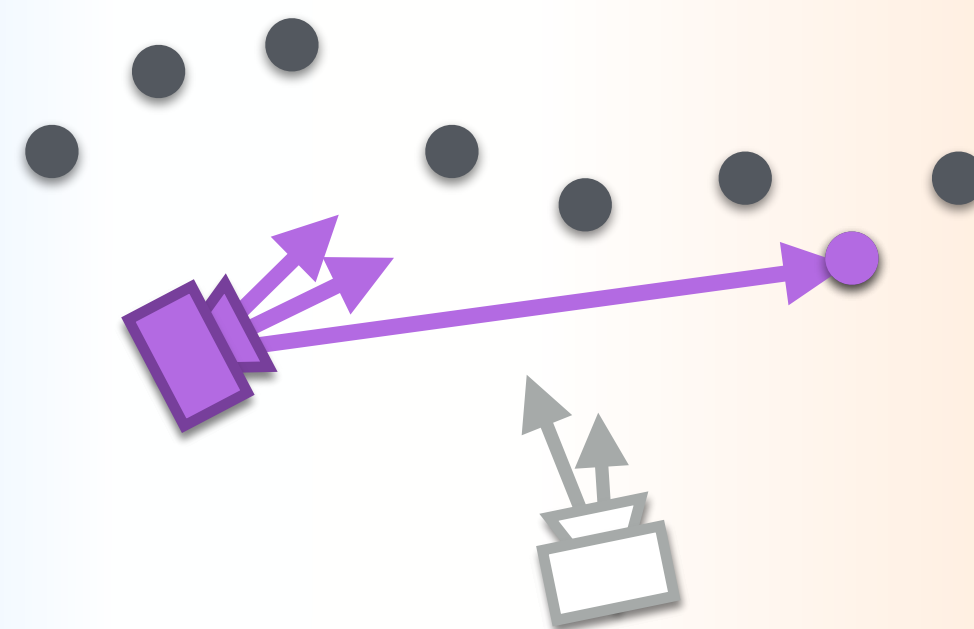


image-based representation

pose triangulation

[Zhang & Kosecka, 3DPVT'06]

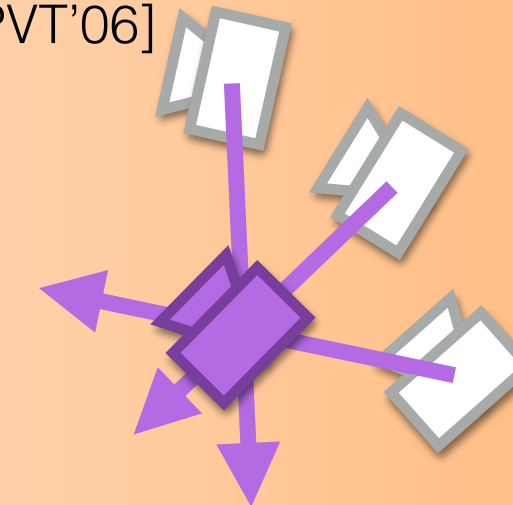
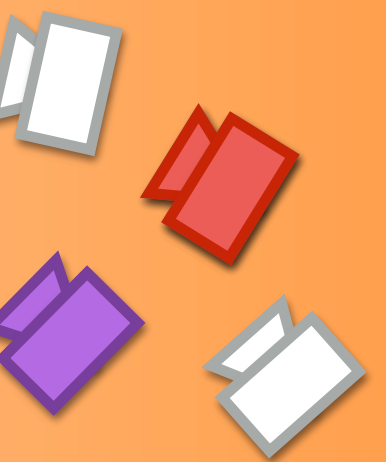
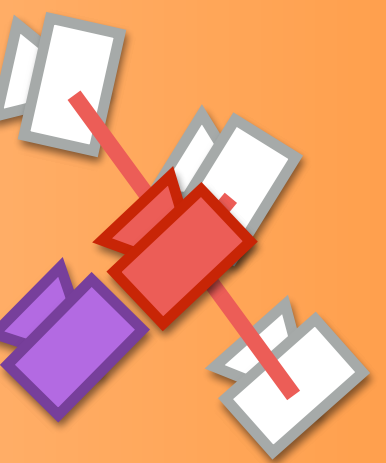


image retrieval



pose interpolation

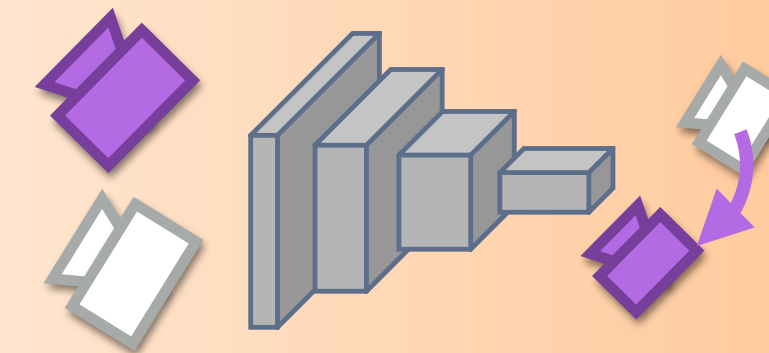


[Torii et al., ICCV'11]

semi-generalized
relative pose / homography

[Zheng & Wu, ICCV'15]

[Bhayani et al., ICCV'21]

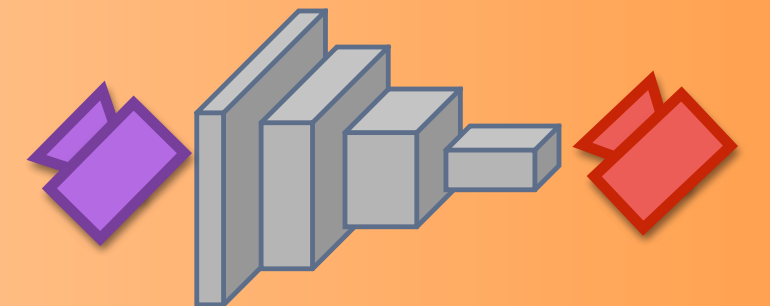


relative pose regression

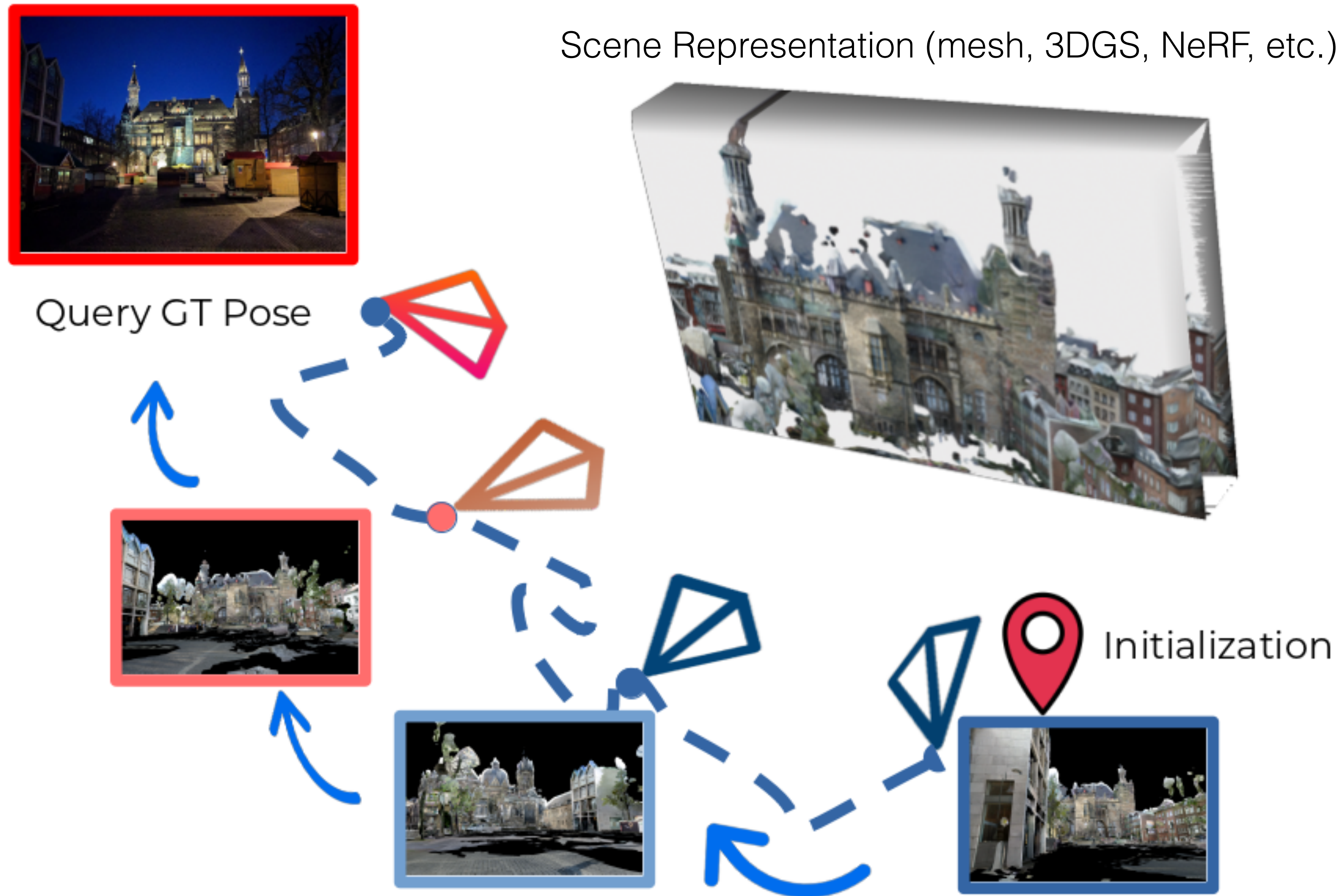
[Laskar et al., CVPRW'17] [Balntas et al., ECCV'18]

absolute pose regression

[Kendall et al., ICCV'15]



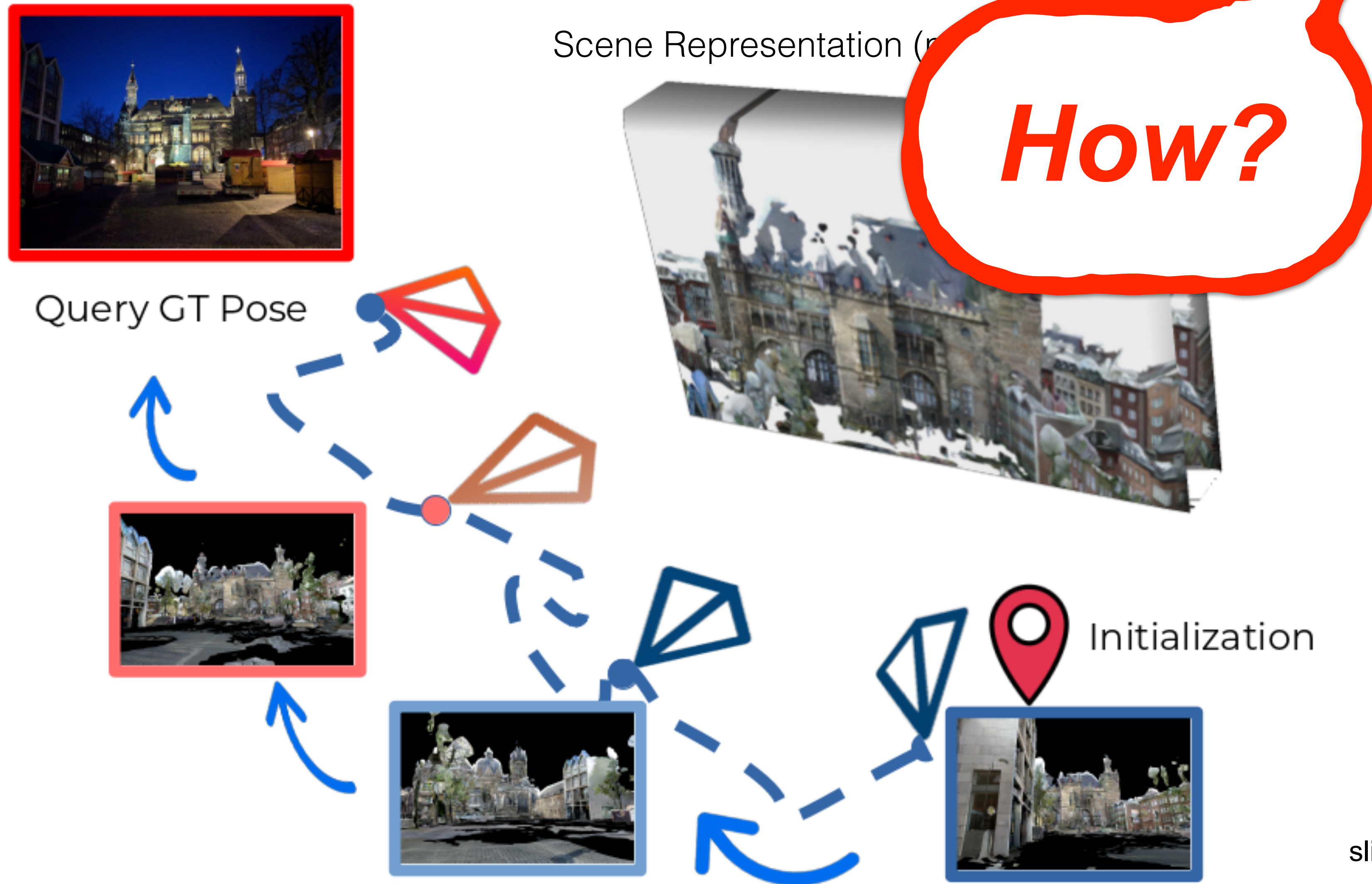
Visual Localization via Render&Compare



slide credit: Gabriele Trivigno

[Gabriele Trivigno, Carlo Masone, Barbara Caputo, Sattler, The Unreasonable Effectiveness of Pre-Trained Features for Camera Pose Refinement, CVPR 2024] (highlight)

Visual Localization via Render&Compare



slide credit: Gabriele Trivigno

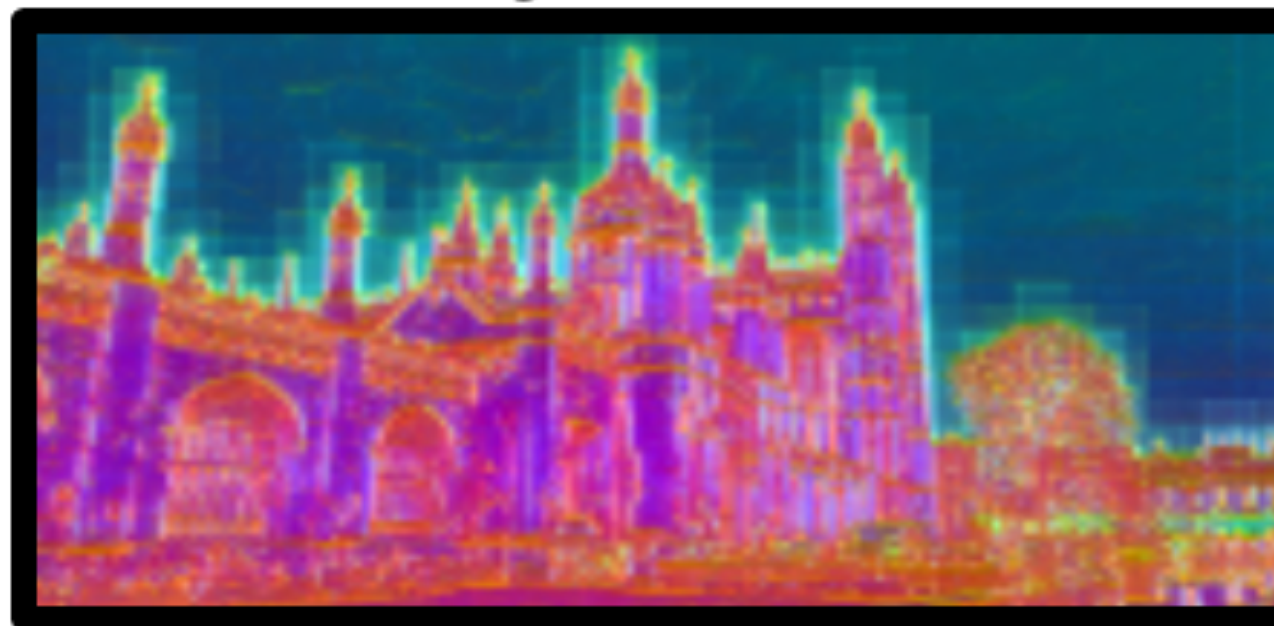
[Gabriele Trivigno, Carlo Masone, Barbara Caputo, Sattler, The Unreasonable Effectiveness of Pre-Trained Features for Camera Pose Refinement, CVPR 2024] (highlight)

Jointly Training Representation and Features

Query Image

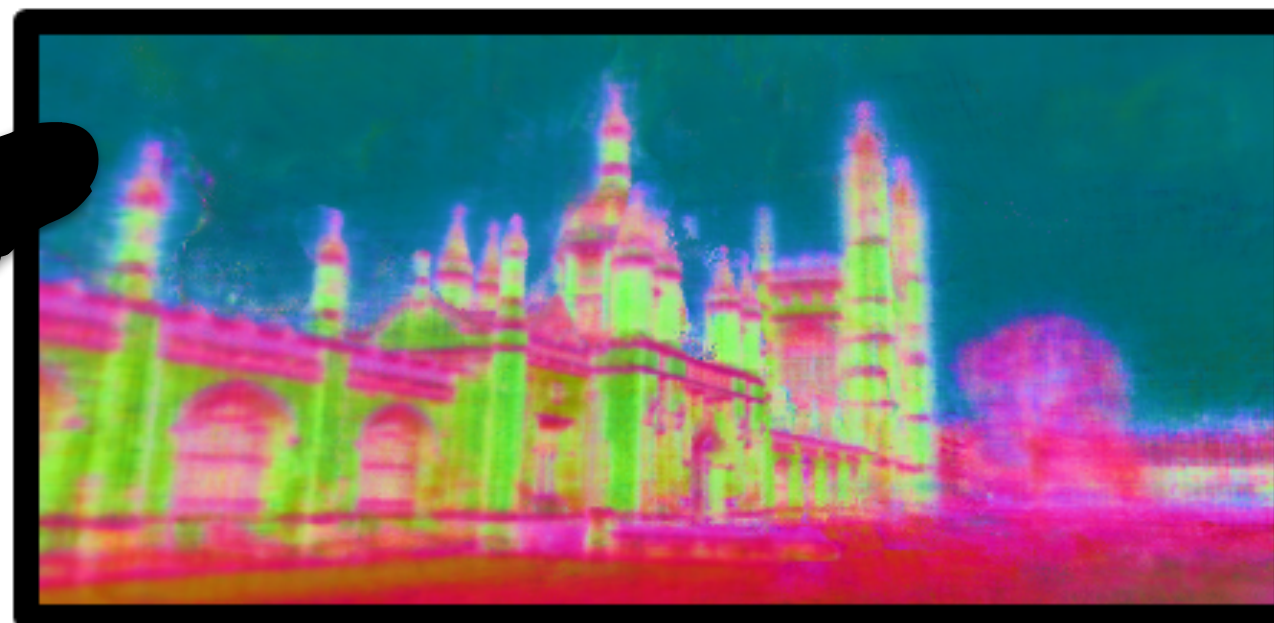


Query Features

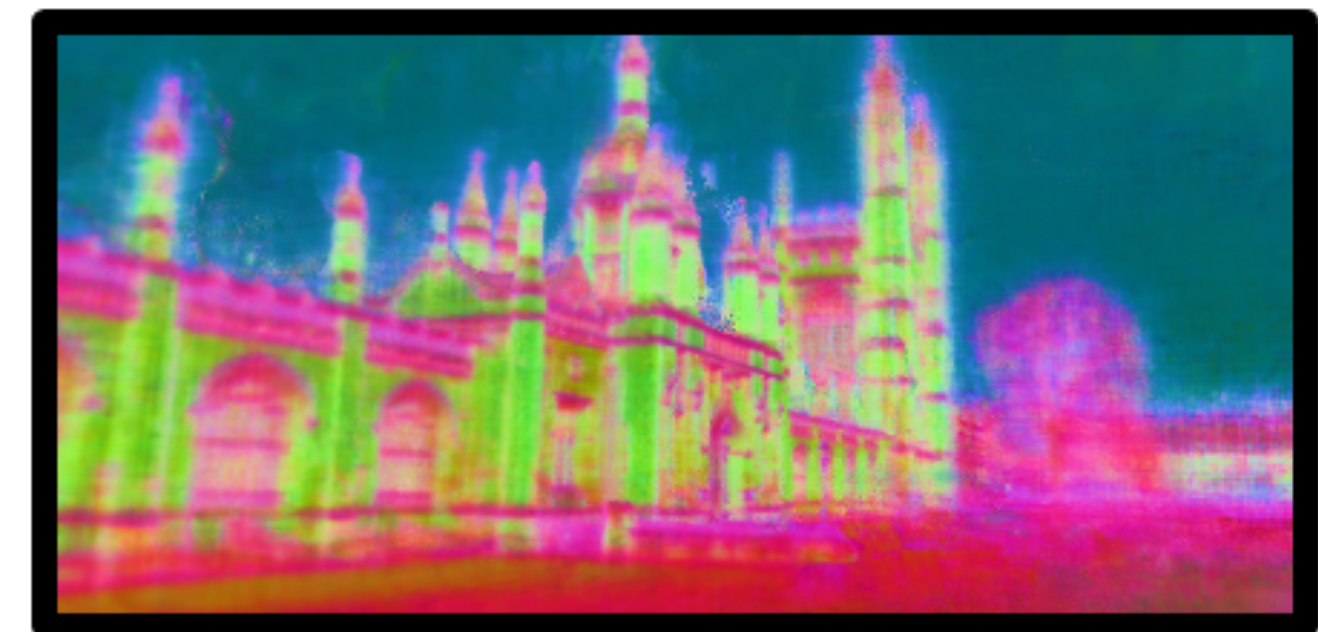


Encoder

Makes sense if
representation not
already given



Rendered
Features Init



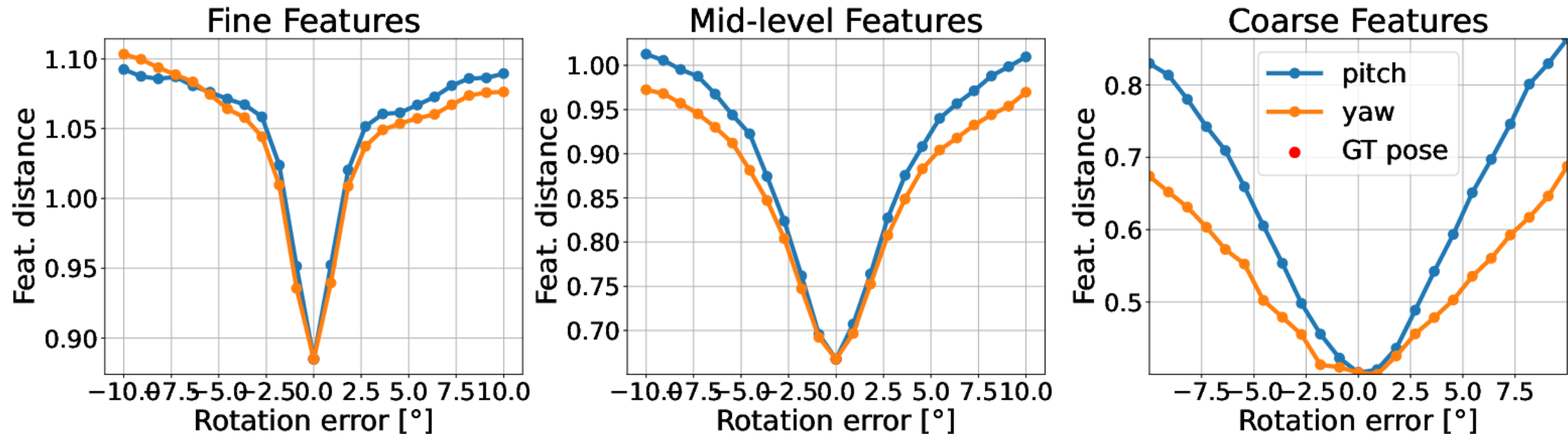
Rendered
Features

Alignment

[Maxime Pietrantoni, Gabriela Csurka, Martin Humenberger, Sattler, Self-supervised learning of Neural implicit Feature Fields for Camera Pose Refinement, 3DV 2024]

The Unreasonable Effectiveness of Pre-Trained Features

- Dense deep features are known to be good estimators of **perceptual similarity**
- This property can be exploited to measure **pose similarity** as well
- **Feature depth** is correlated with the **sensitivity** → hierarchical scheme



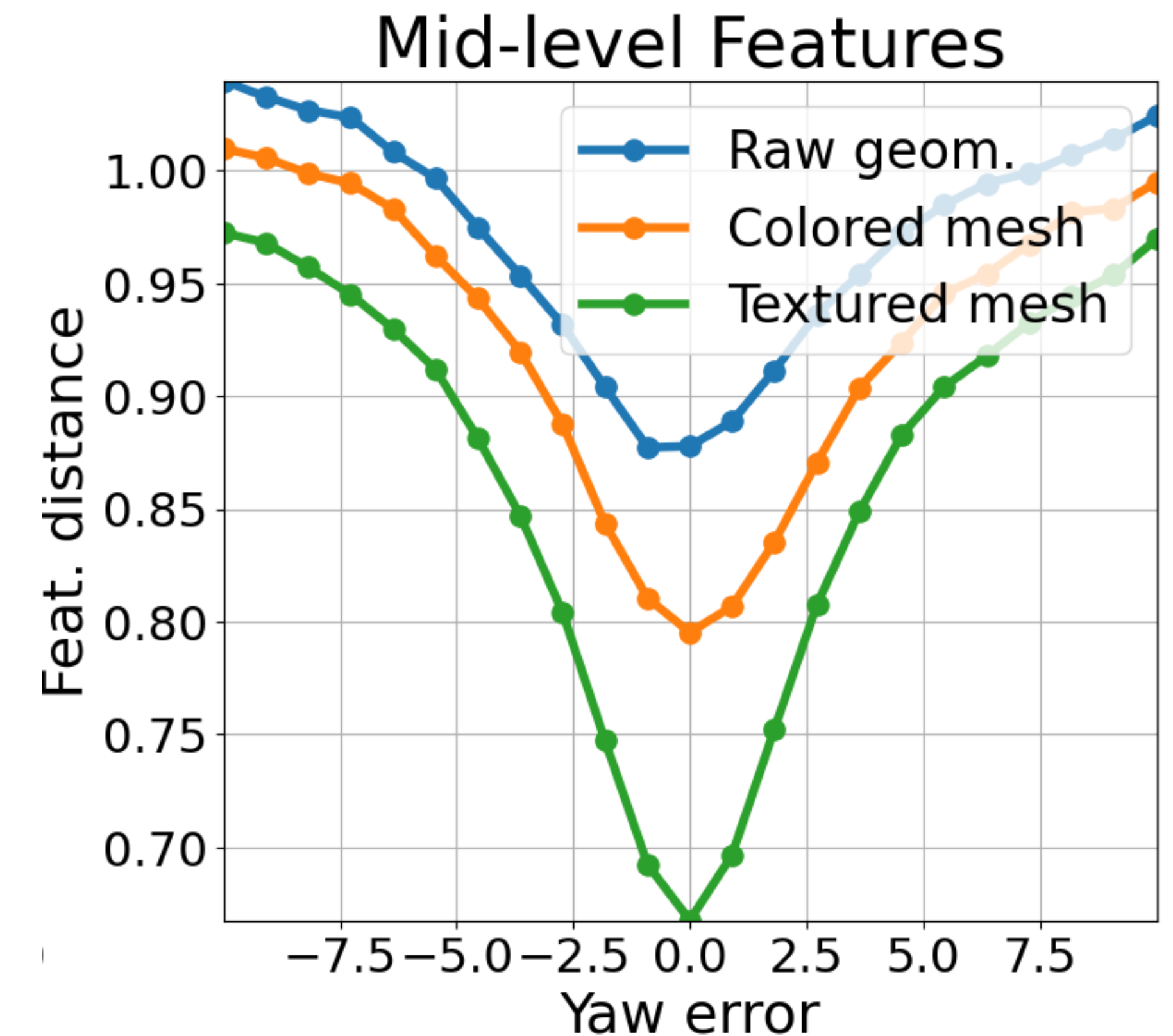
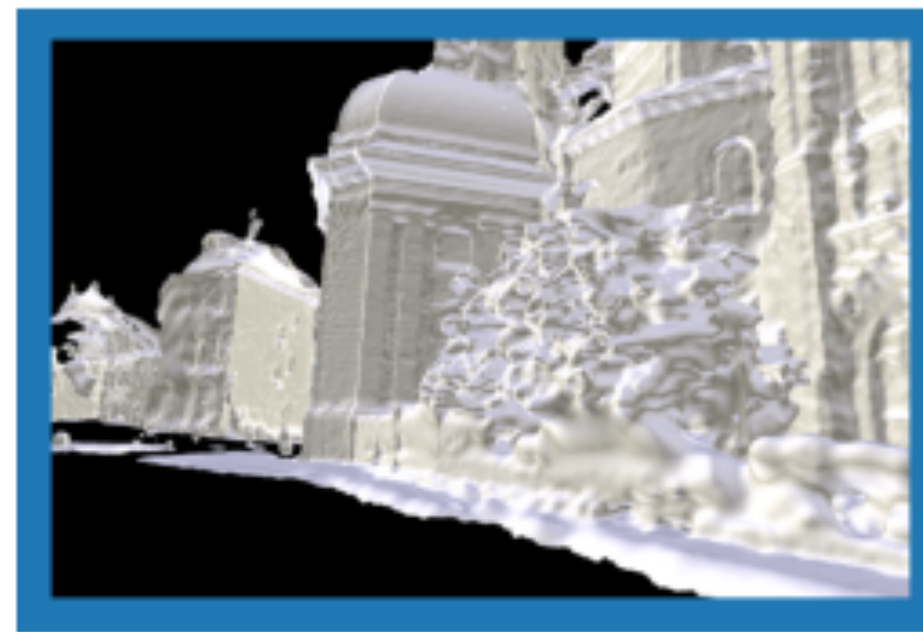
slide credit: Gabriele Trivigno

[Gabriele Trivigno, Carlo Masone, Barbara Caputo, Sattler, The Unreasonable Effectiveness of Pre-Trained Features for Camera Pose Refinement, CVPR 2024]

The Unreasonable Effectiveness of Pre-Trained Features

Dense deep features are quite robust to domain changes

Query



slide credit: Gabriele Trivigno

[Gabriele Trivigno, Carlo Masone, Barbara Caputo, Sattler, The Unreasonable Effectiveness of Pre-Trained Features for Camera Pose Refinement, CVPR 2024]

Pose Refinement based on Rendering

Advantages:

- ✓ Improves good initial poses
- ✓ Can handle poor geometry (depth not directly used)

Disadvantages:

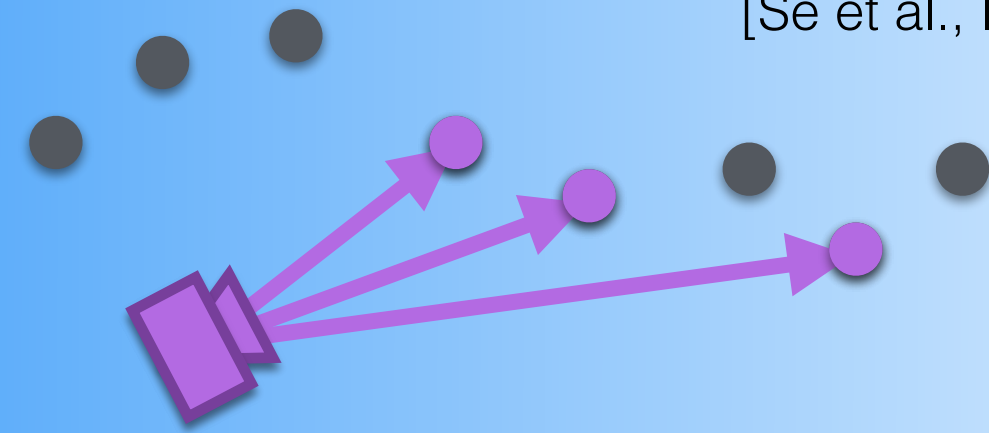
- ✗ Accuracy depends on initialization
- ✗ Basis of convergence limited
- ✗ Can be quite slow

Visual Localization - A Taxonomy

3D structure-based representation

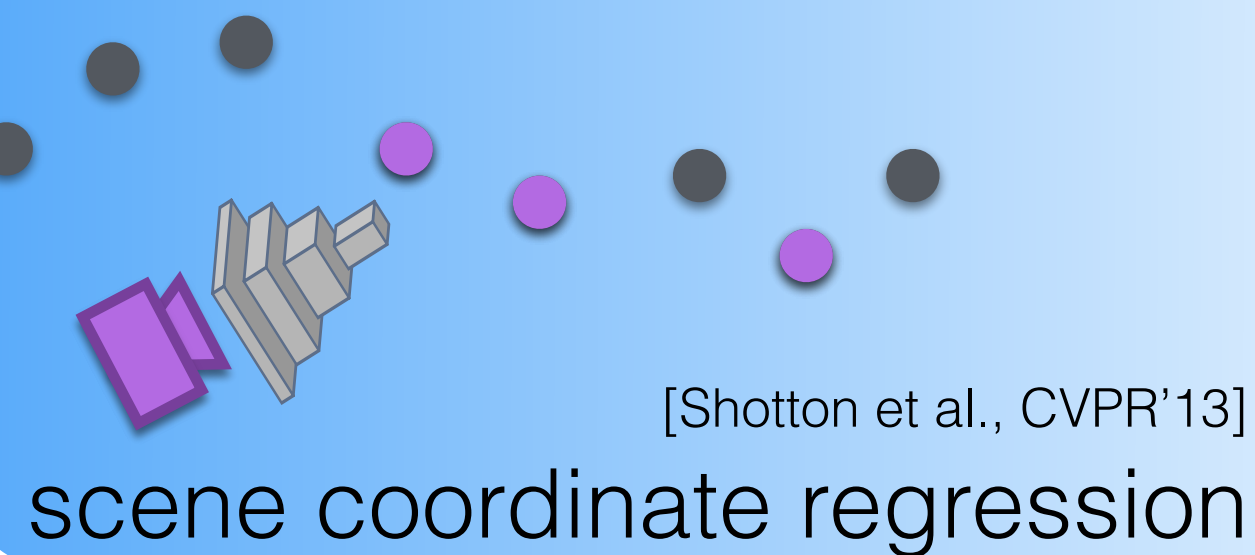
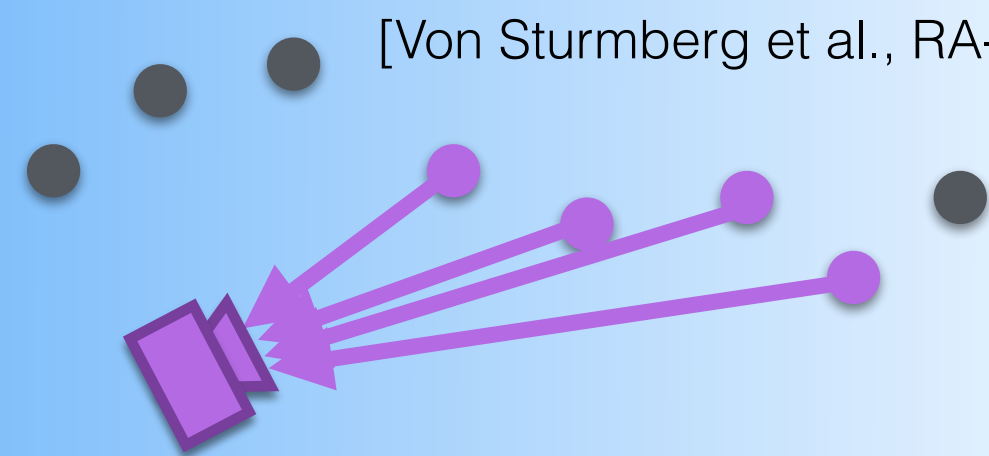
feature-based localization

[Se et al., IROS'02]



pose refinement

[Von Sturmberg et al., RA-L'20]



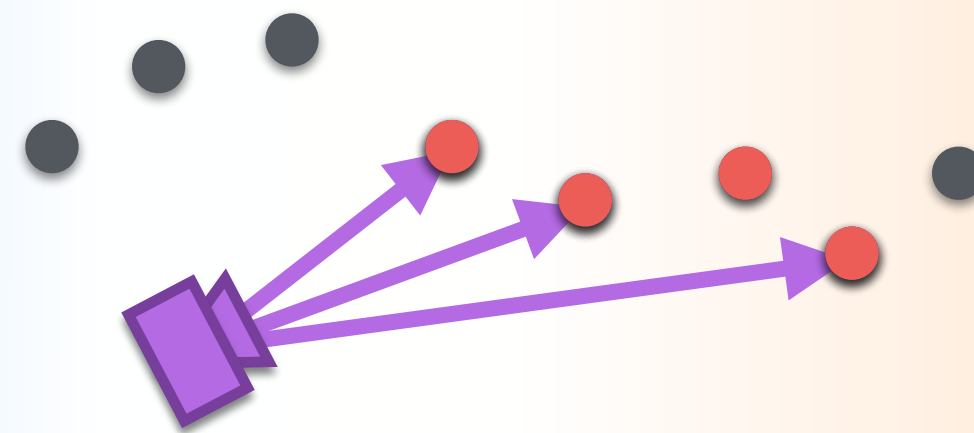
[Shotton et al., CVPR'13]

scene coordinate regression

hybrid

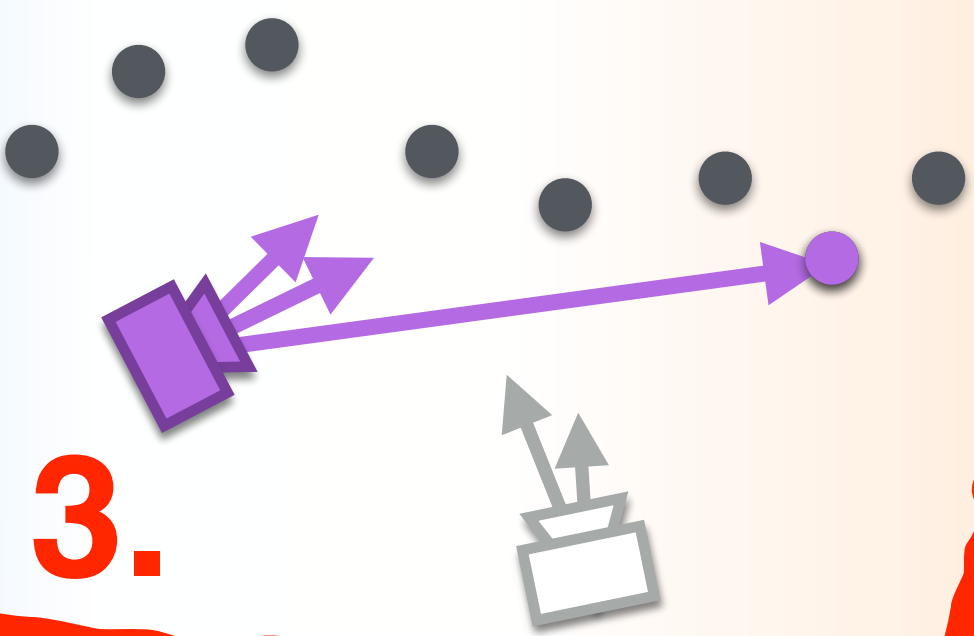
hierarchical localization

[Irschara et al., CVPR'09]



hybrid pose estimation

[Camposeco et al., ICCV'21]



3.

image-based representation

pose triangulation

[Zhang & Kosecka, 3DPVT'06]

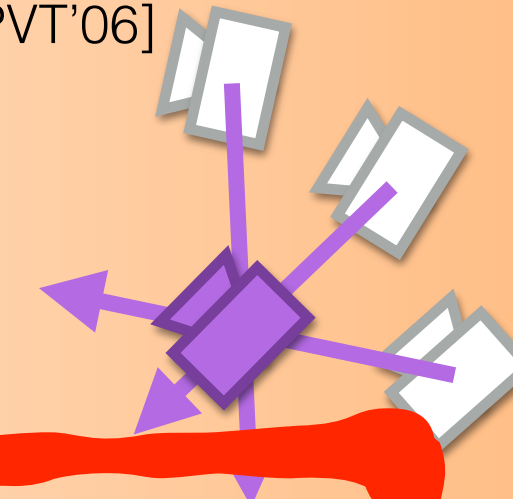
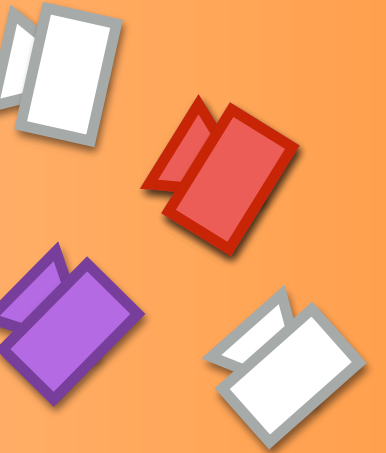
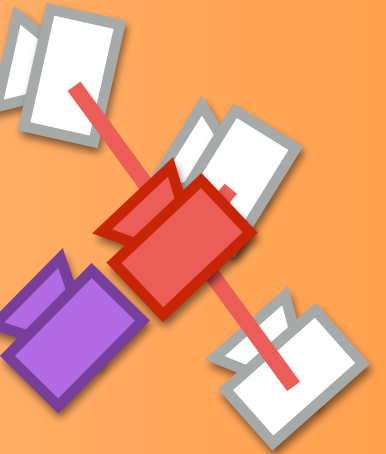


image retrieval



pose interpolation

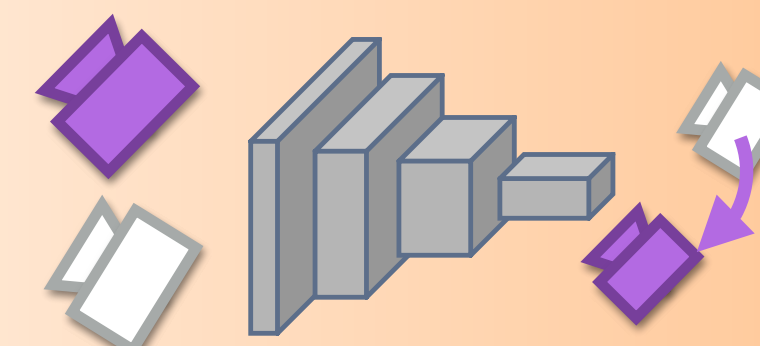


[Torii et al., ICCVW'11]

semi-generalized relative pose / homography

[Zheng & Wu, ICCV'15]

[Bhayani et al., ICCV'21]

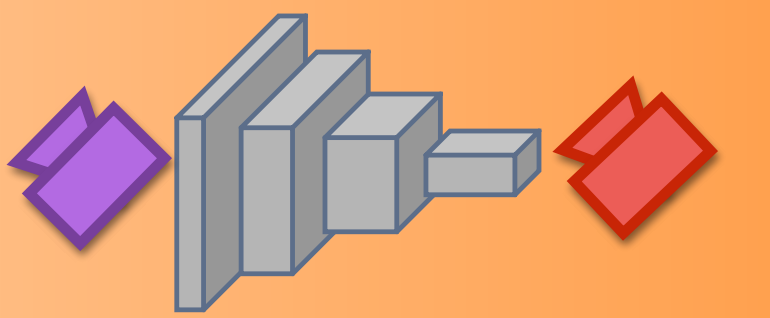


relative pose regression

[Laskar et al., CVPRW'17] [Balntas et al., ECCV'18]

absolute pose regression

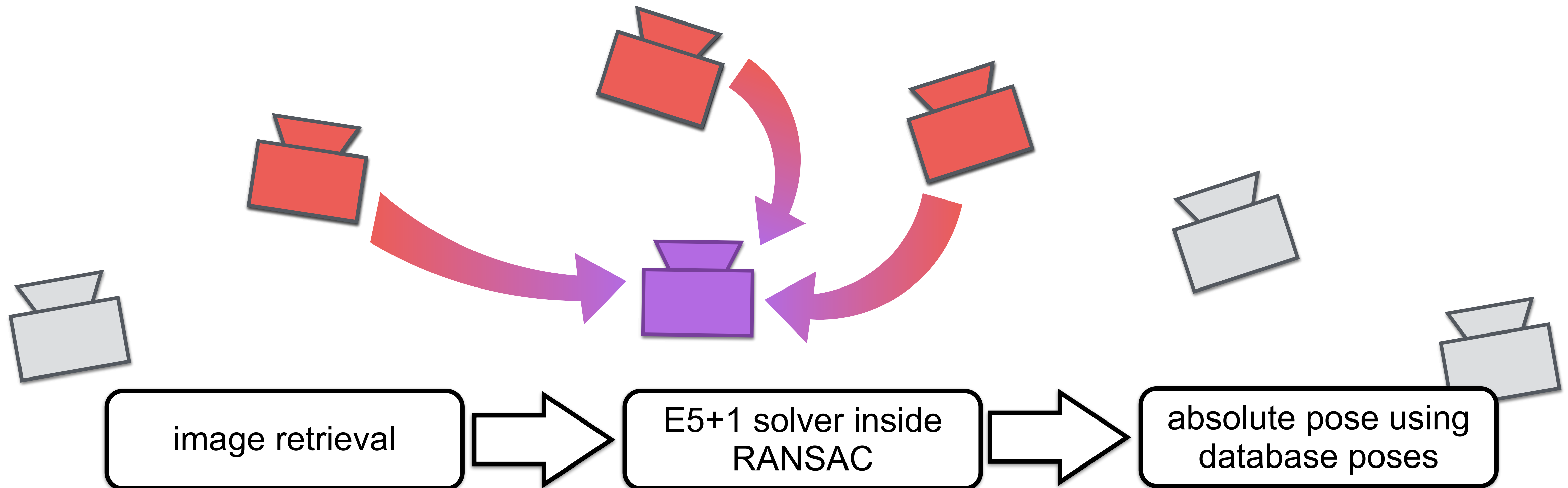
[Kendall et al., ICCV'15]



Structure-Less Visual Localization

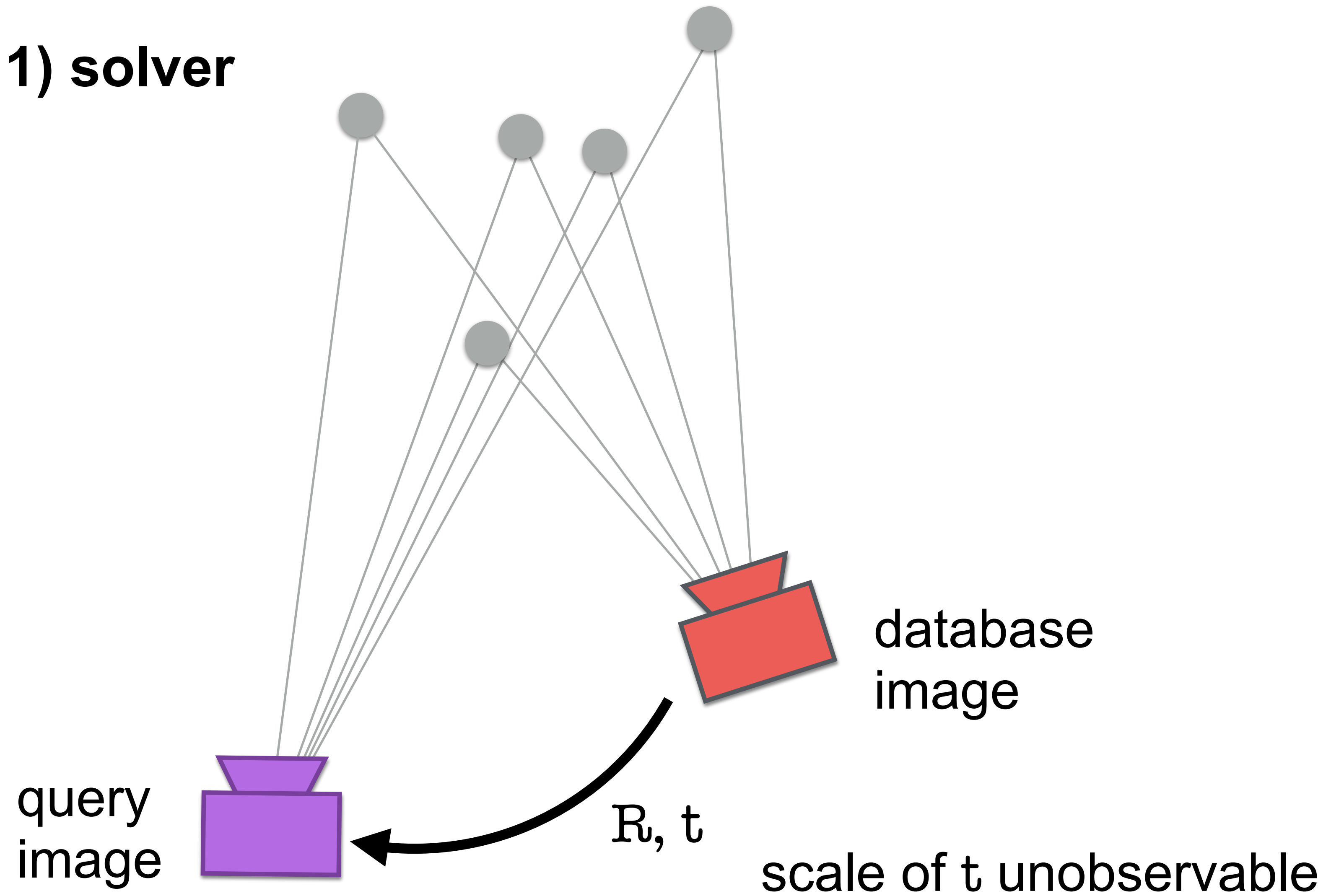
Scene representation: images with known poses, **no 3D points**

- ✓ Easy to update: just add / remove image and pose from database
- ✓ Extract features on the fly, easy to use new feature type



Semi-Generalized Relative Pose Estimation

E5+1 (Essential Matrix + 1) solver

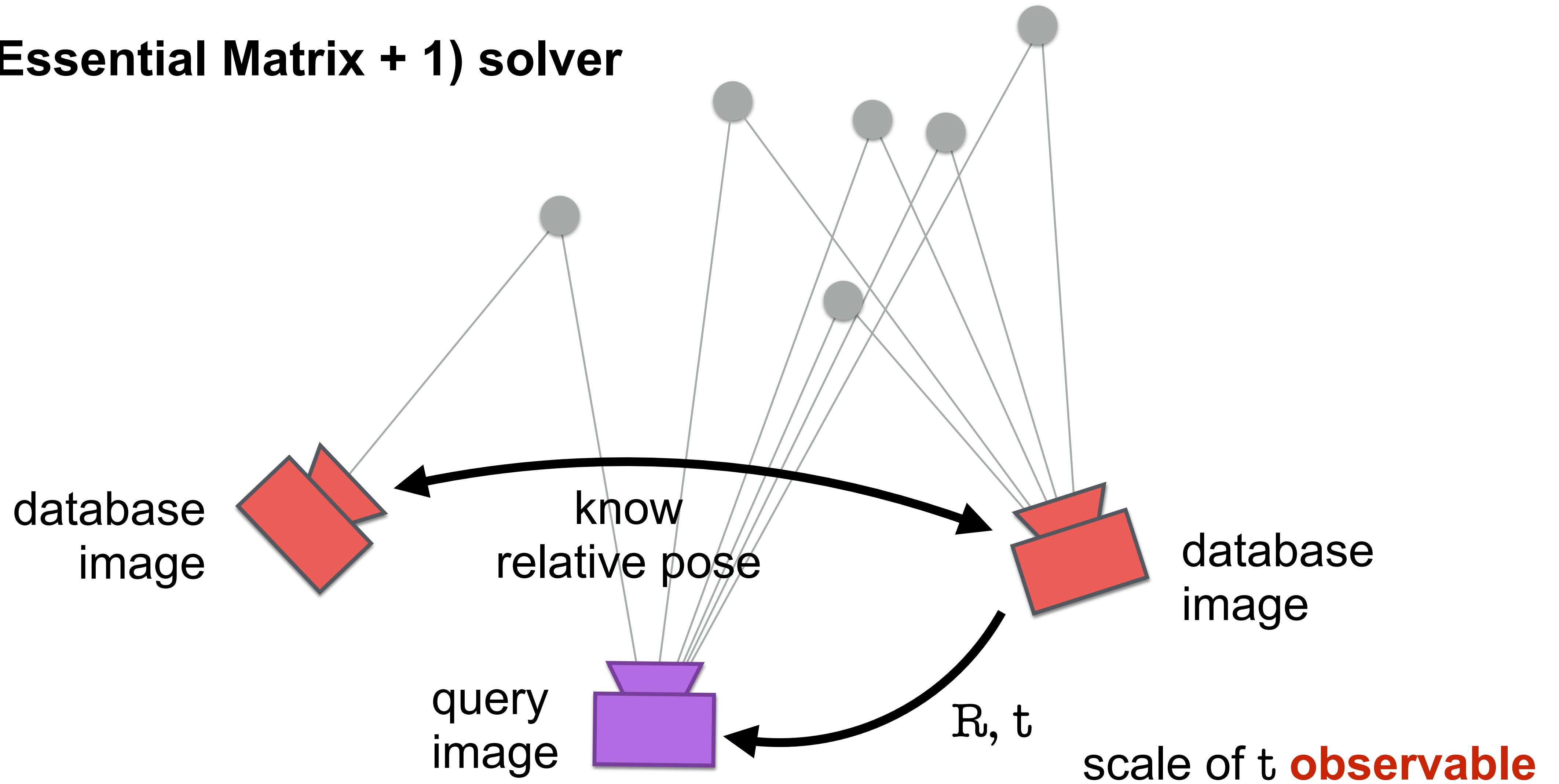


[Zheng, Wu, Structure from Motion Using Structure-less Resection, ICCV 2015]

[Bhayani, Sattler, Barath, Beliansky, Heikkila, Kukelova, Calibrated and Partially Calibrated Semi-Generalized Homographies, ICCV 2021]

Semi-Generalized Relative Pose Estimation

E5+1 (Essential Matrix + 1) solver

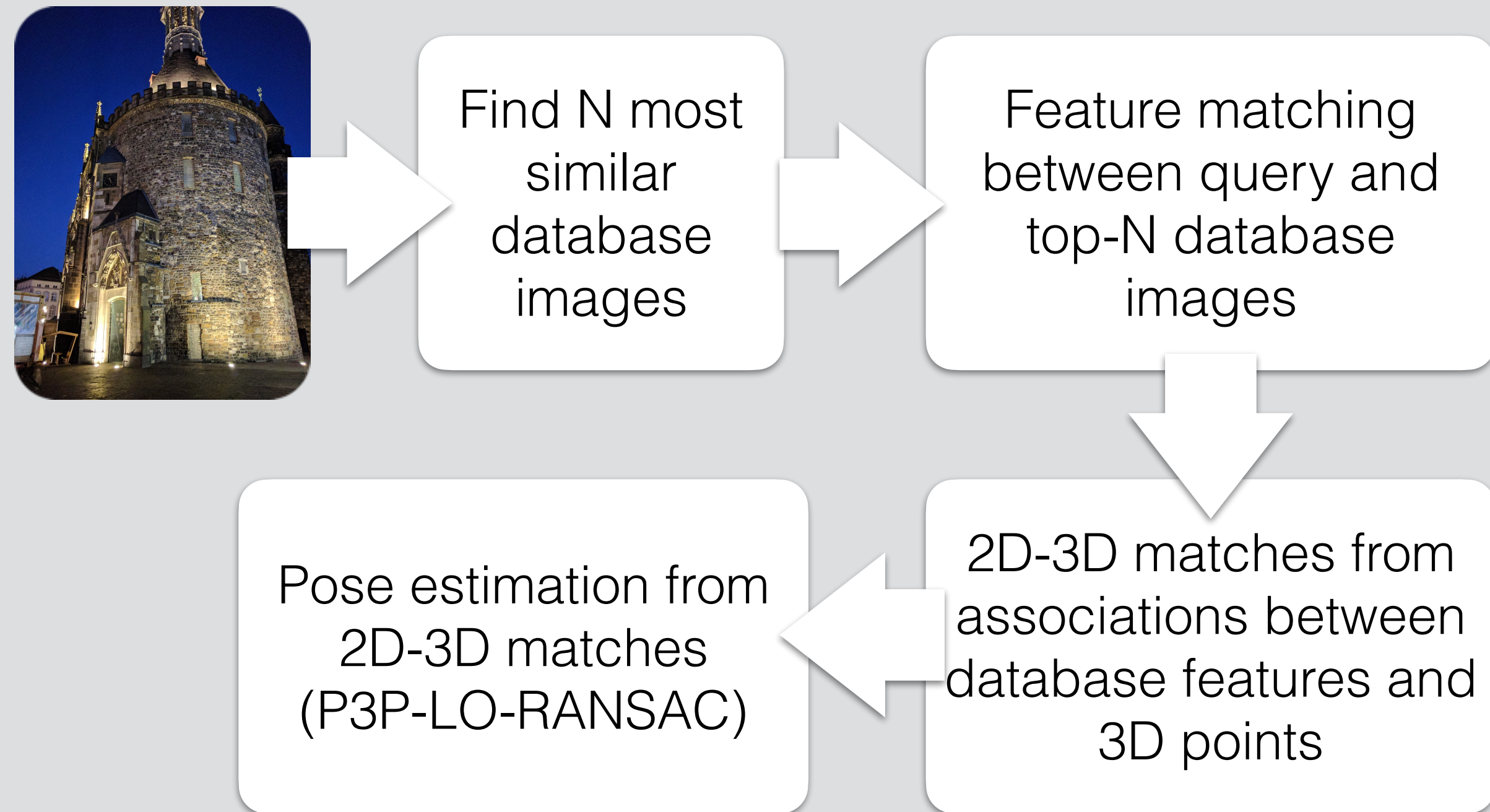


[Zheng, Wu, Structure from Motion Using Structure-less Resection, ICCV 2015]

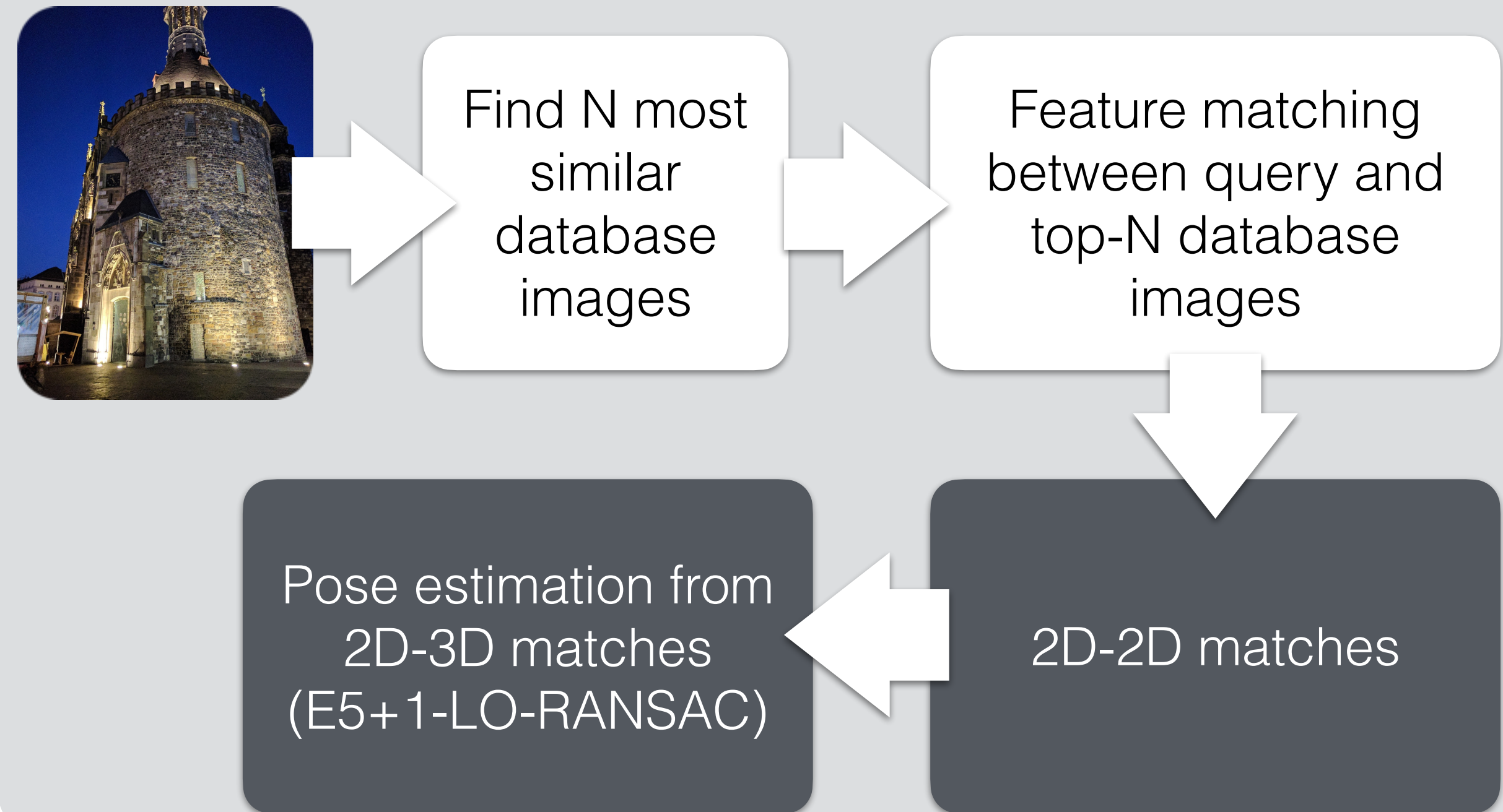
[Bhayani, Sattler, Barath, Beliansky, Heikkila, Kukelova, Calibrated and Partially Calibrated Semi-Generalized Homographies, ICCV 2021]

Structure-Based vs. Structure-Less Localization

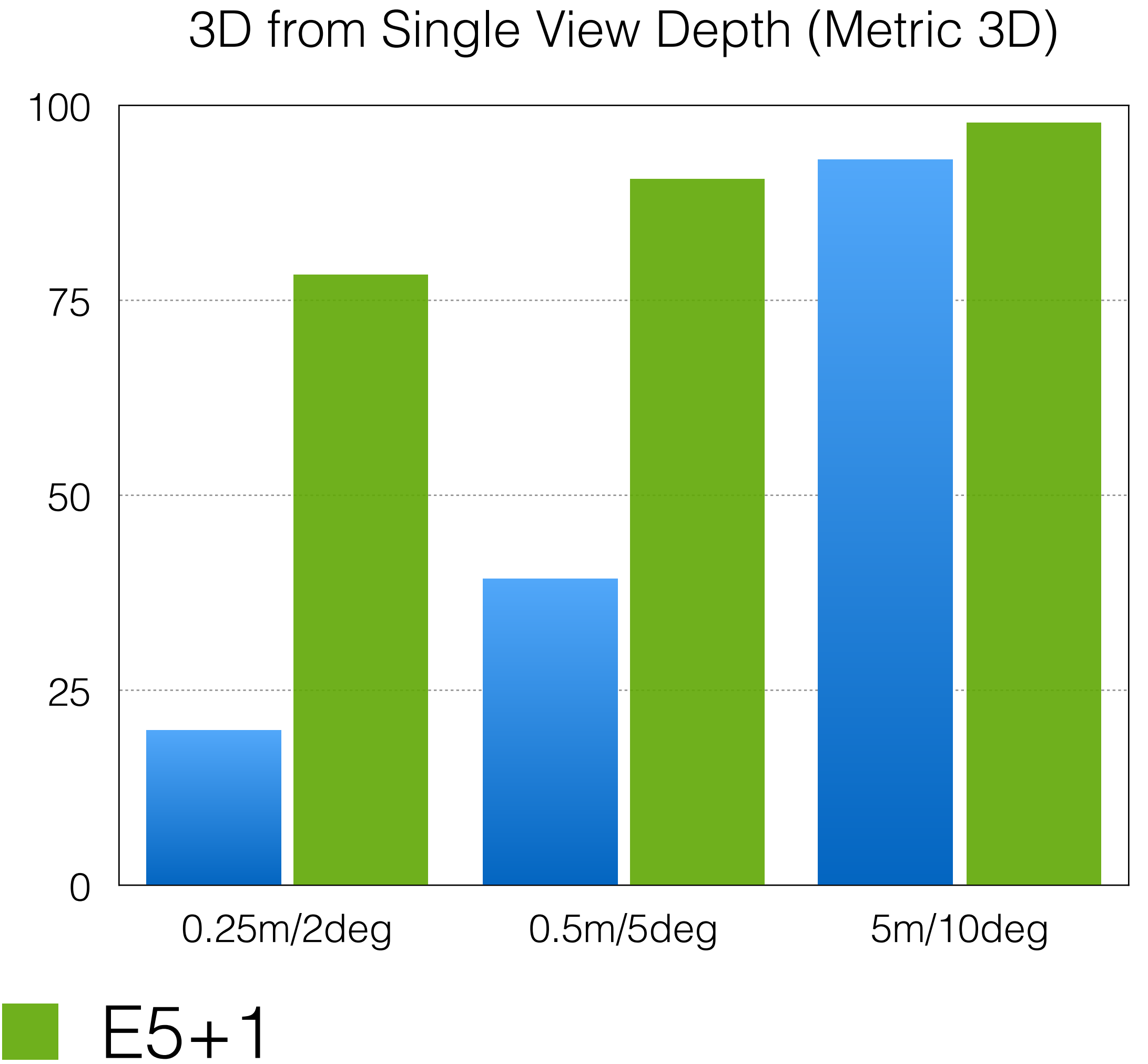
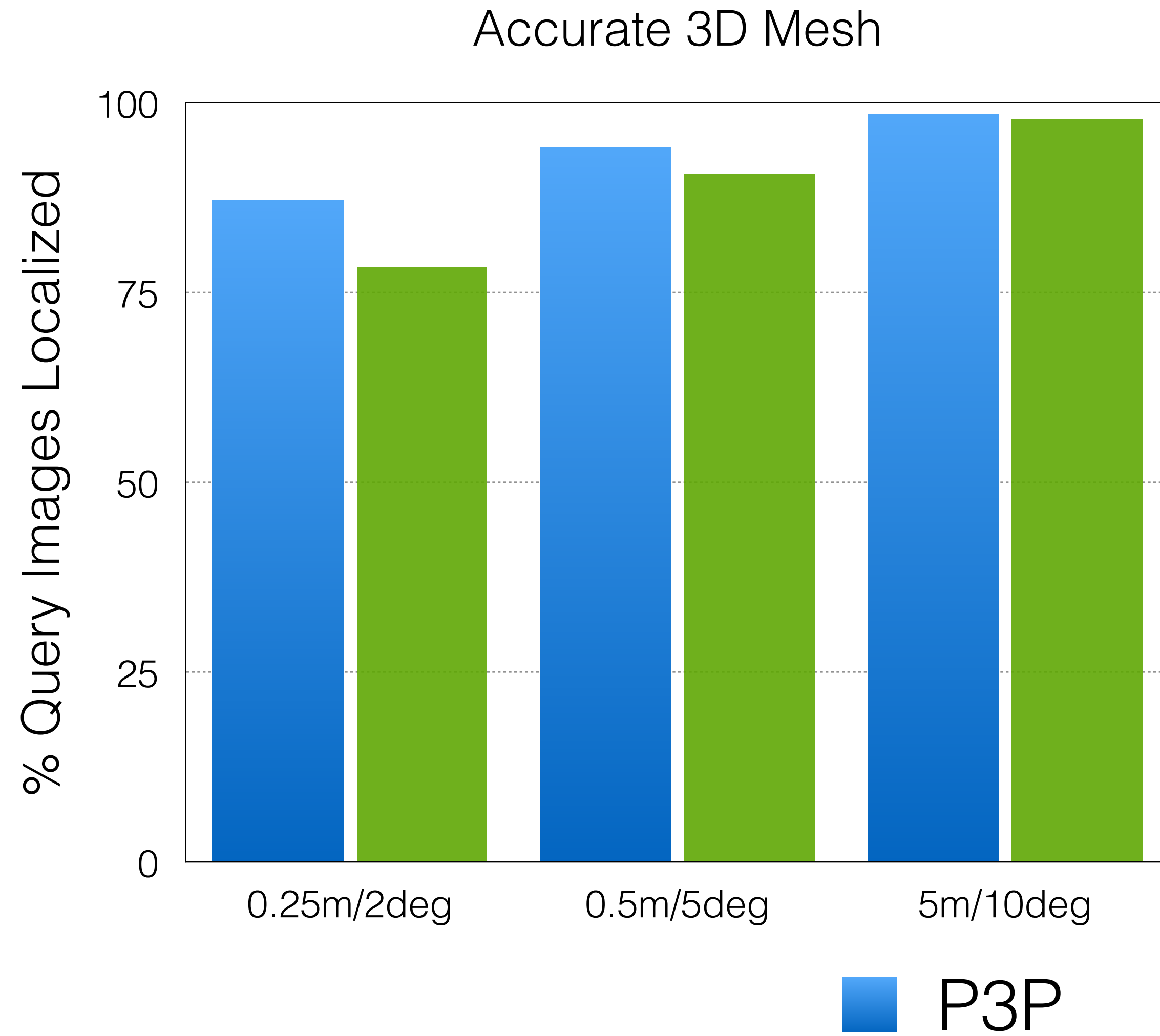
Structure-Based Localization



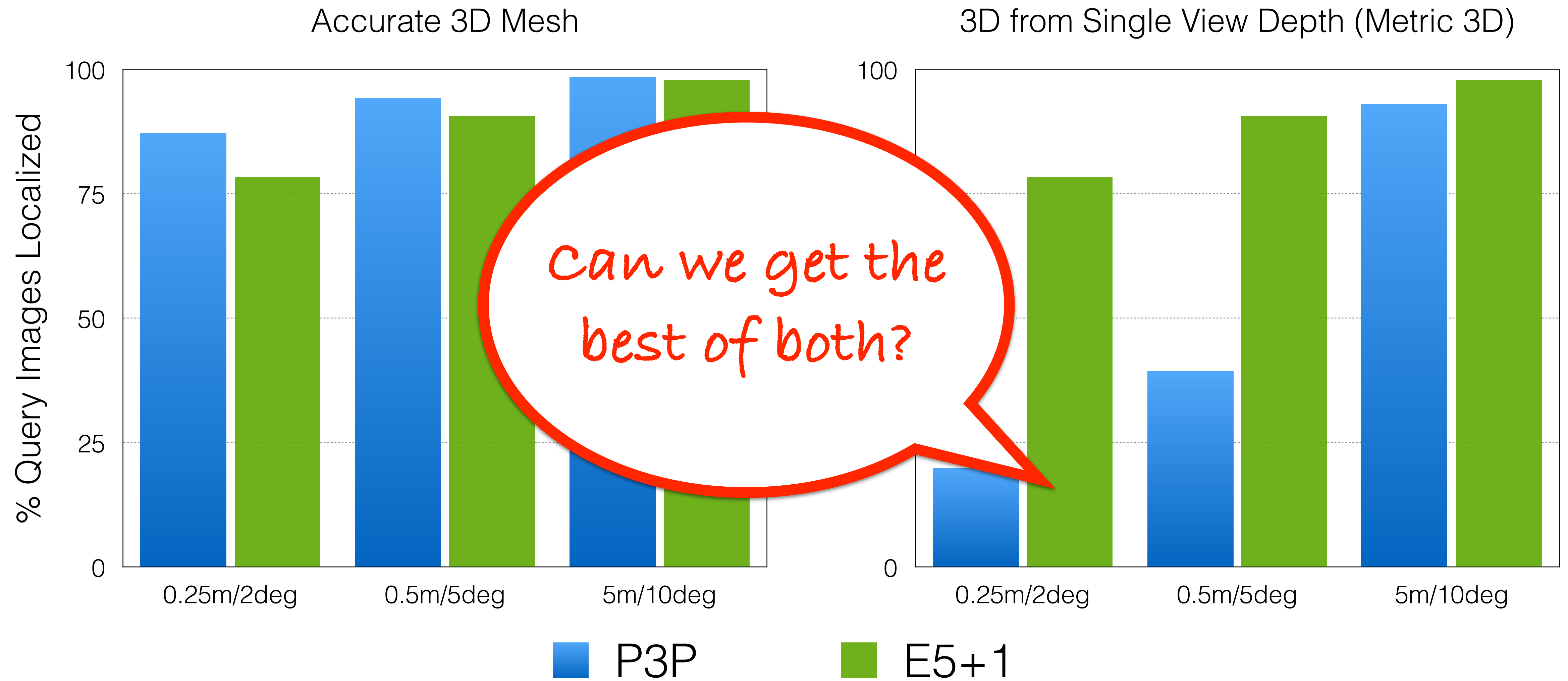
Structure-Less Localization



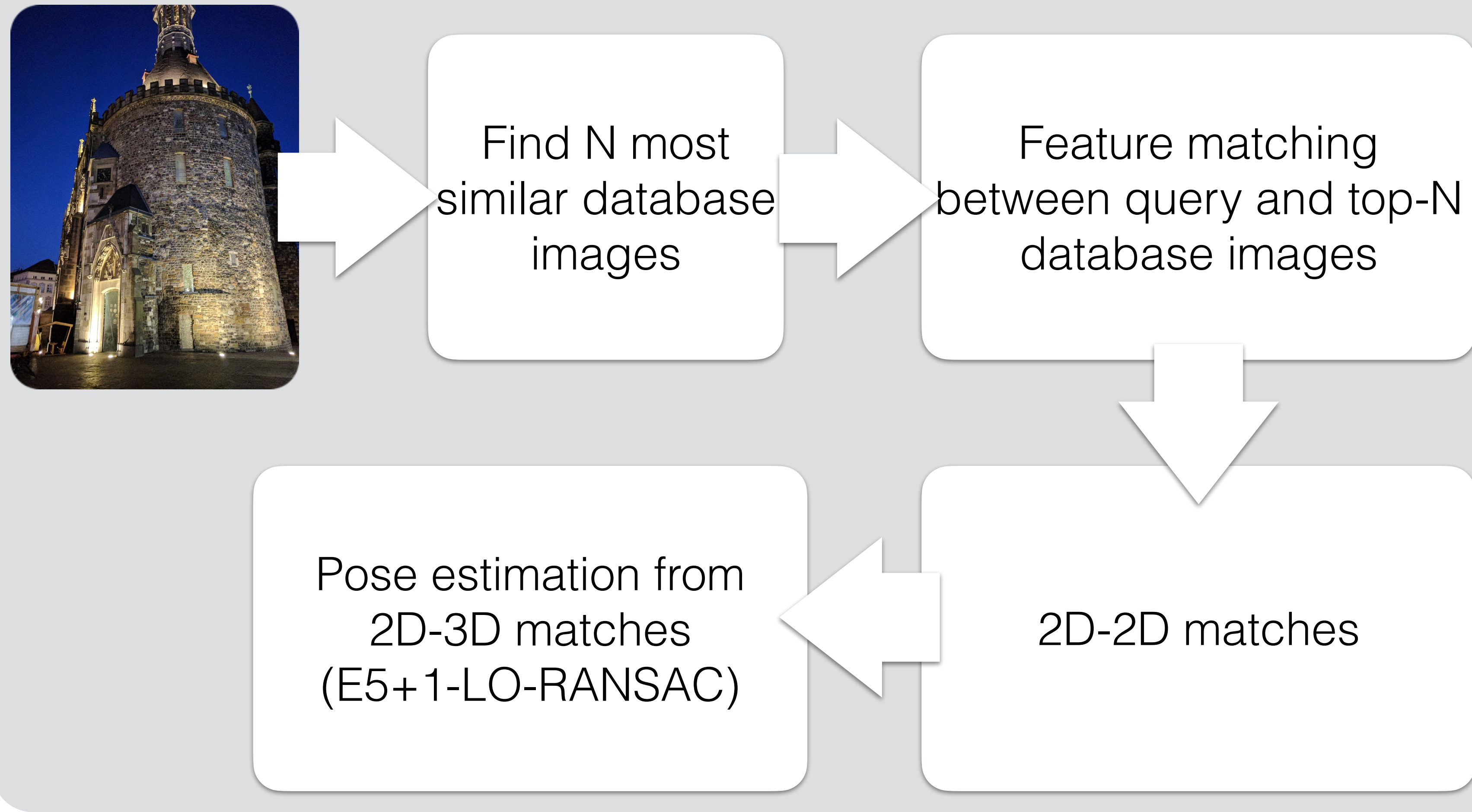
Structure-Based vs. Structure-Less Localization



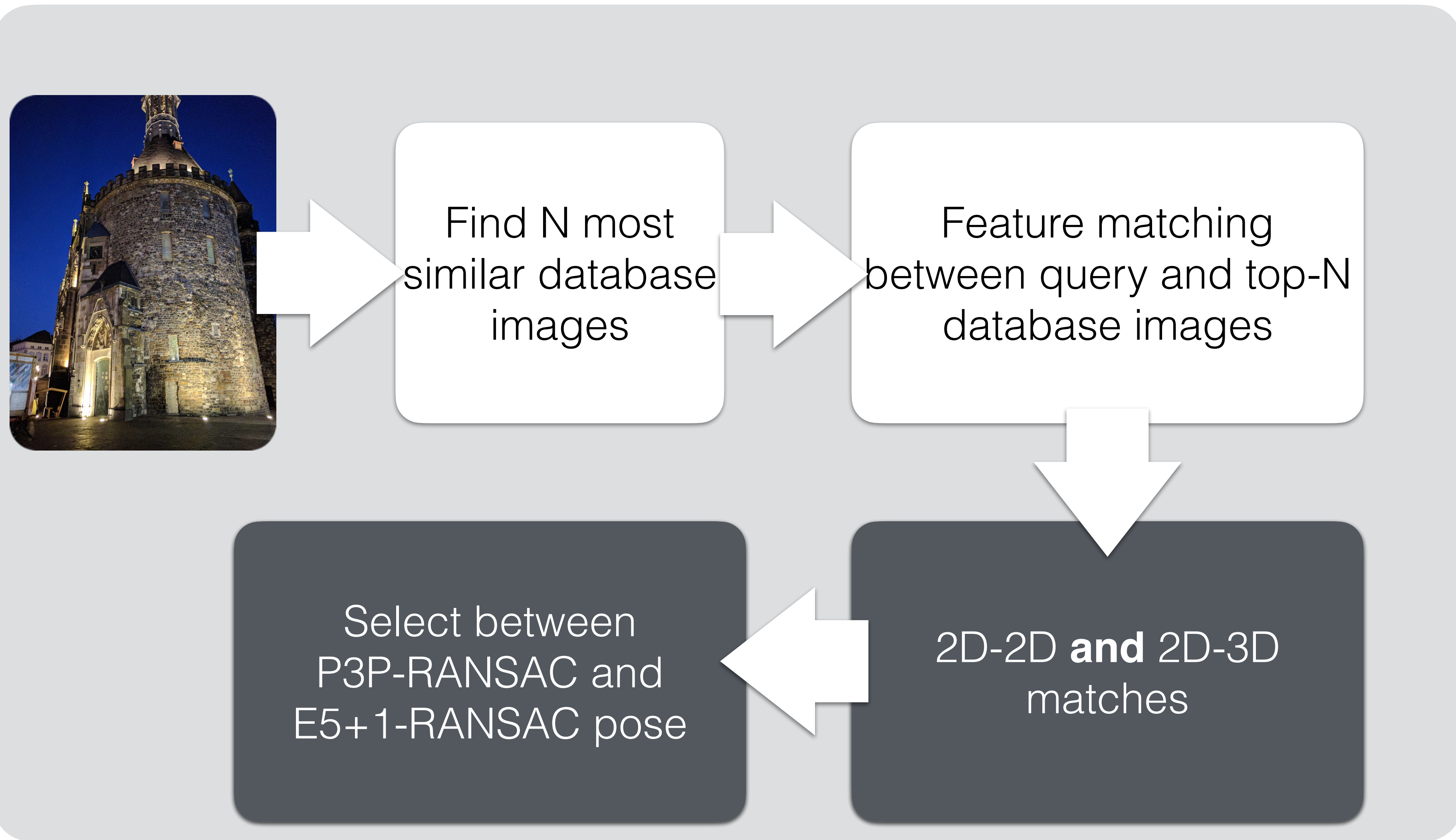
Structure-Based vs. Structure-Less Localization



An Adaptive Strategy



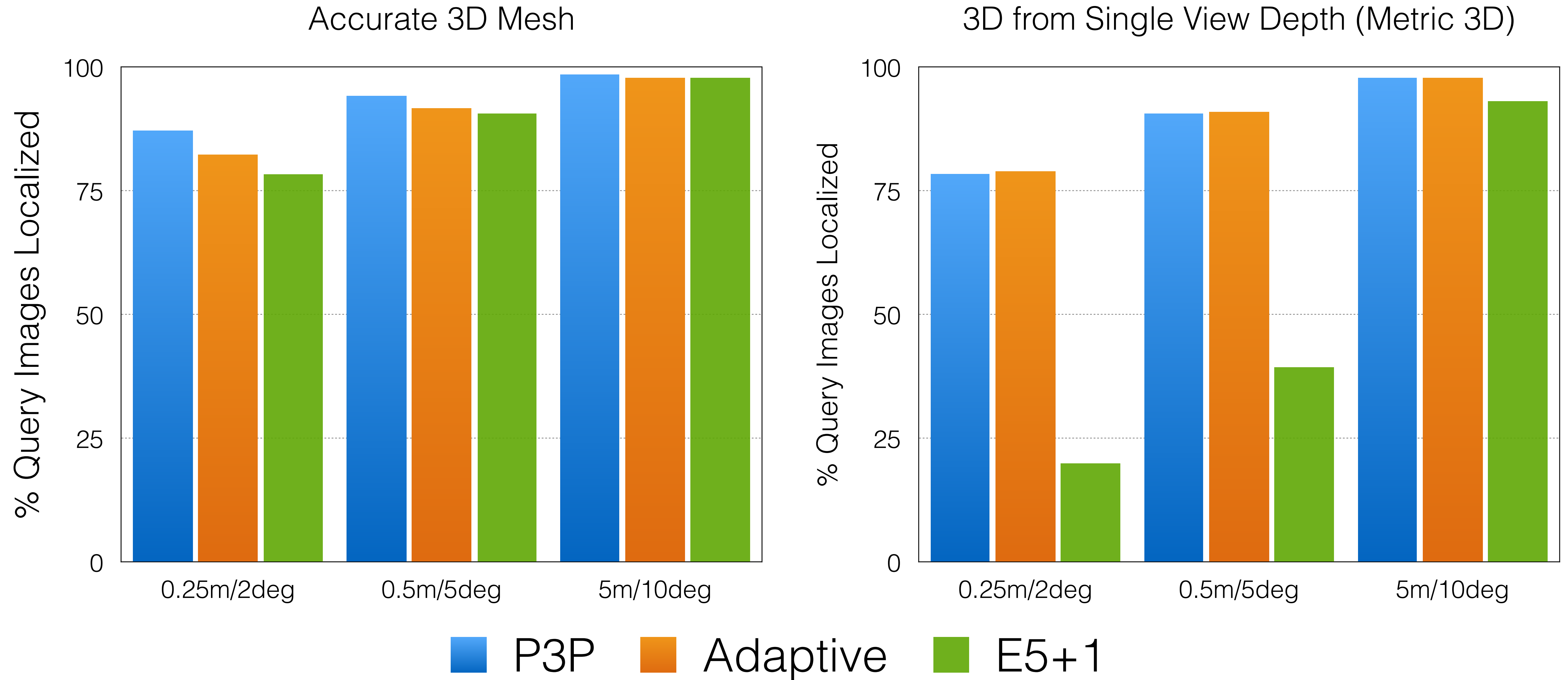
An Adaptive Strategy



[Panek, Sattler, Kukulova, Combining Absolute and Semi-Generalized Relative Poses for Visual Localization, arXiv:2409.14269]

Torsten Sattler

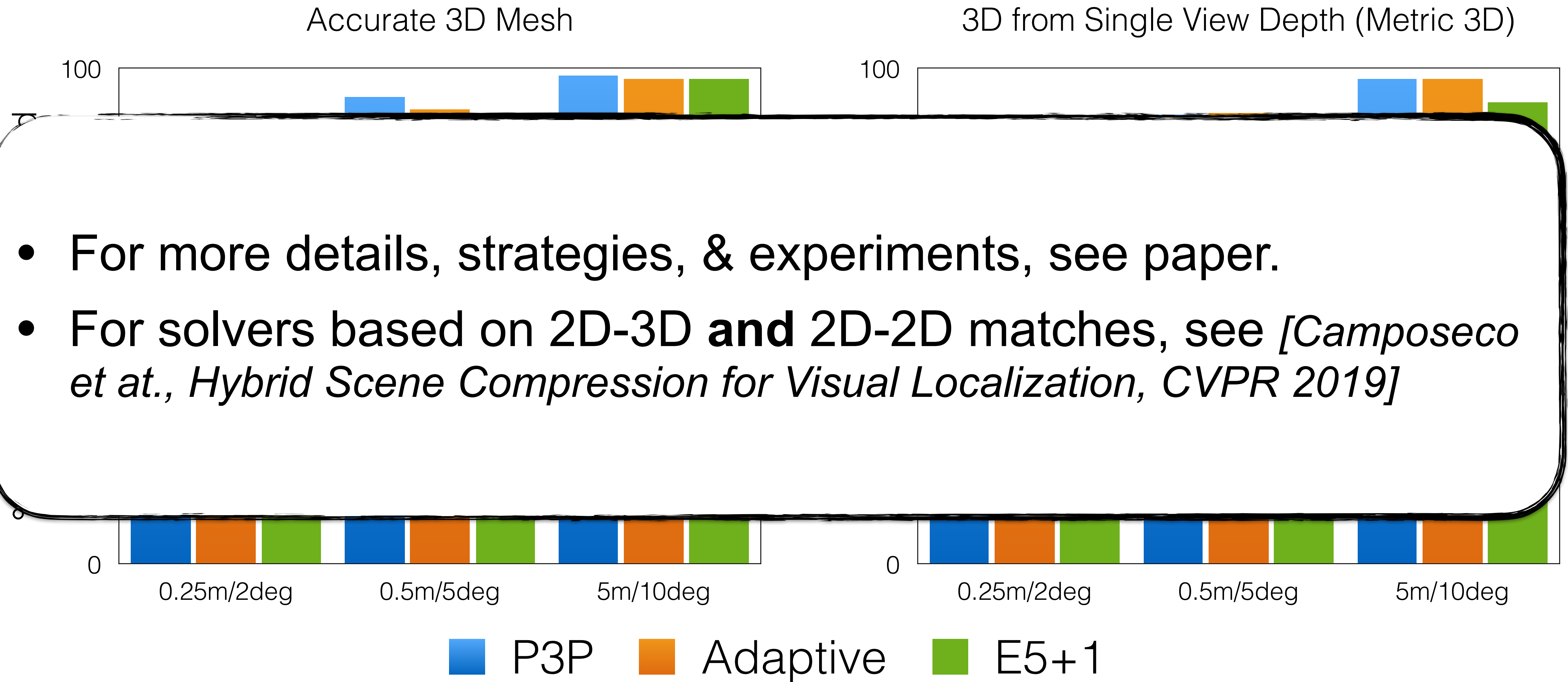
Structure-Less vs. Structure-Based Localization



[Panek, Sattler, Kukelova, Combining Absolute and Semi-Generalized Relative Poses for Visual Localization, arXiv:2409.14269]

Torsten Sattler

Structure-Less vs. Structure-Based Localization

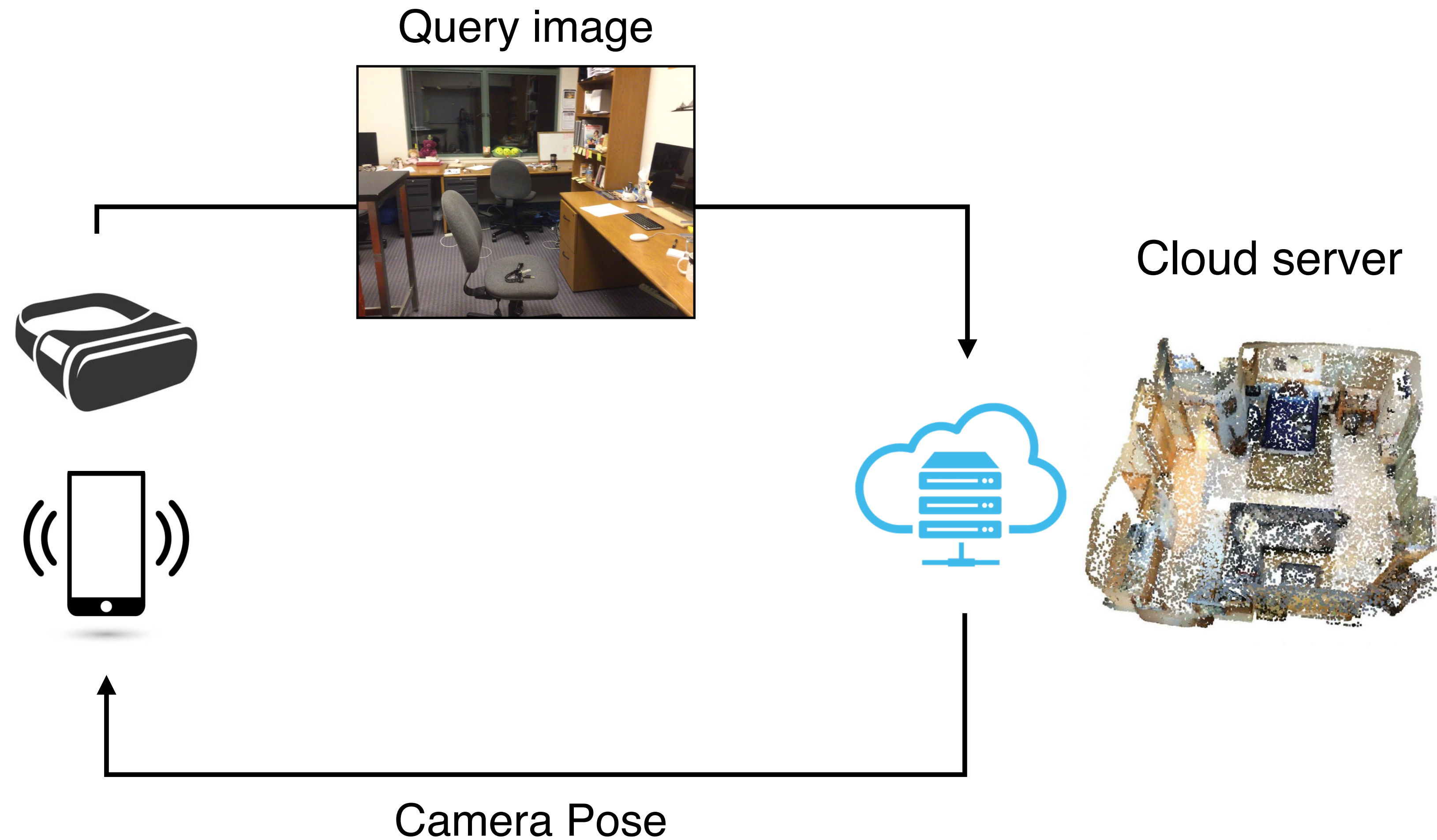


[Panek, Sattler, Kukulova, Combining Absolute and Semi-Generalized Relative Poses for Visual Localization, arXiv:2409.14269]

Torsten Sattler

What About Privacy?

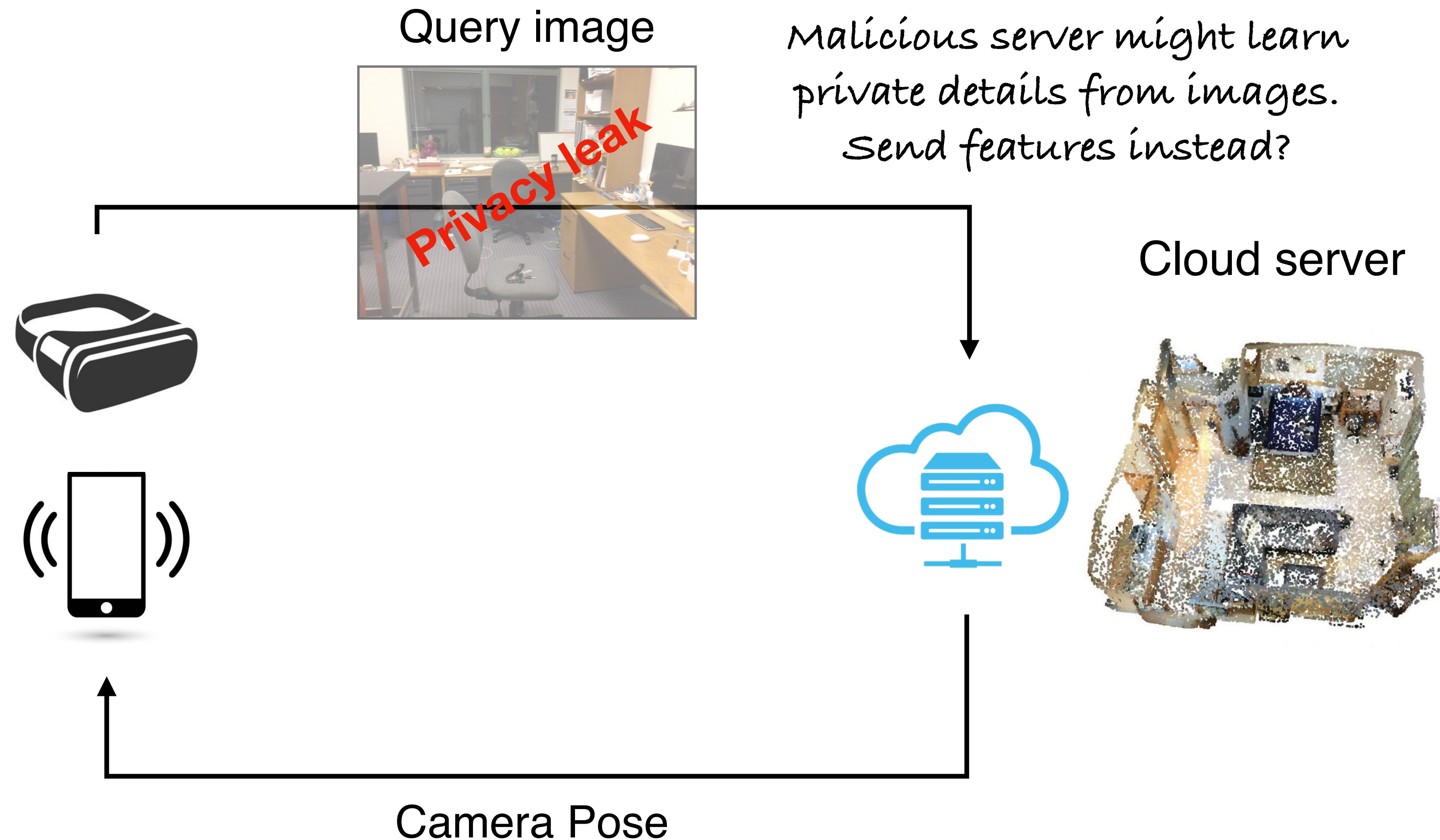
What About Privacy?



slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

What About Privacy?



slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

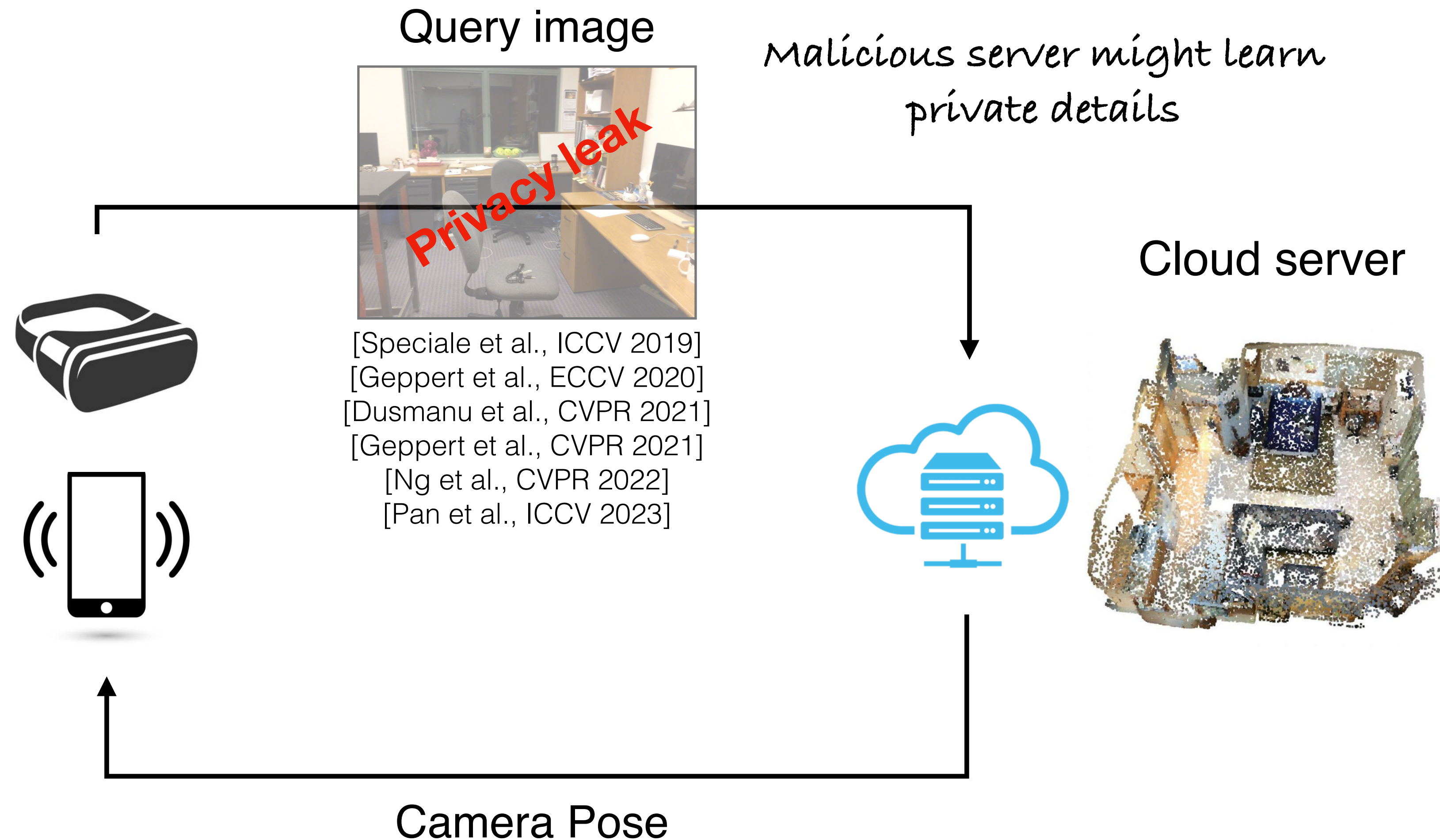
Privacy Issues in Visual Localization

SfM Point Cloud: Scene 1 (NYU)



[Pittaluga, Koppal, Kang, Sinha, Revealing Scenes by Inverting Structure From Motion Reconstructions, CVPR 2019]

What About Privacy?



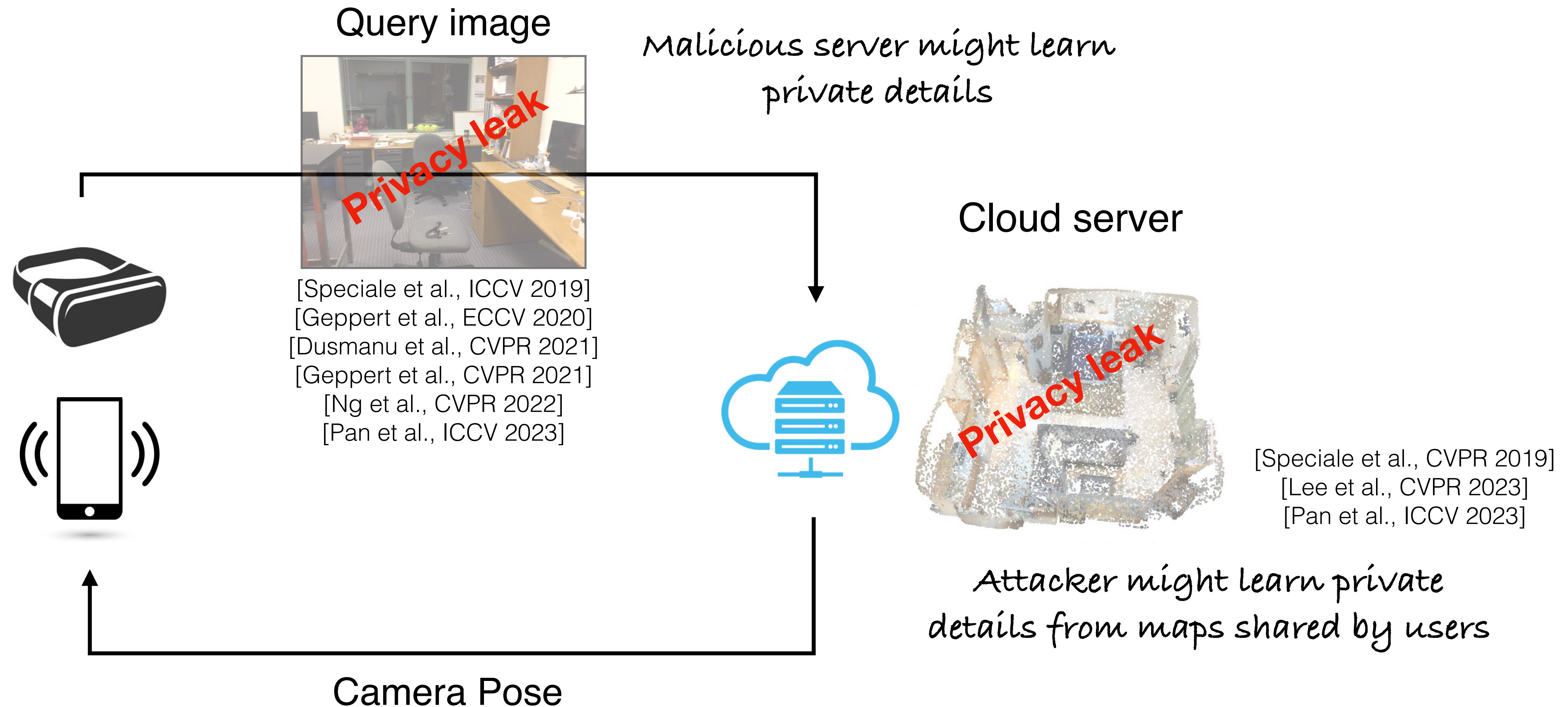
slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

[Chelani, Benbihi, Kahl, Sattler, Kukulova, Obfuscation Based Privacy Preserving Representations are Recoverable Using Neighborhood Information, arXiv:2409.11536]

Torsten Sattler

What About Privacy?



slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

[Chelani, Benbihi, Kahl, Sattler, Kukulova, Obfuscation Based Privacy Preserving Representations are Recoverable Using Neighborhood Information, arXiv:2409.11536]

Torsten Sattler

What About Privacy?

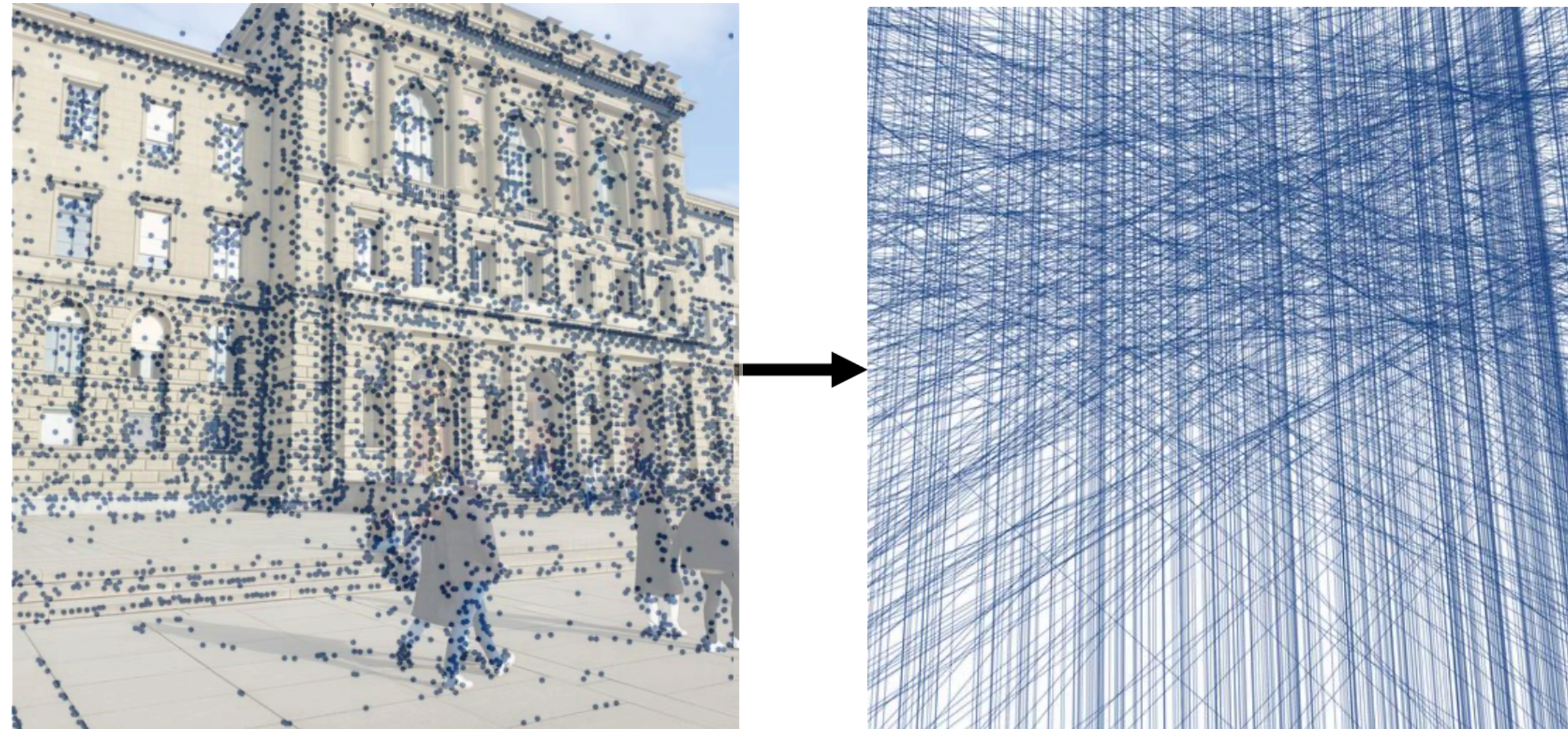


image credit: Marcel Geppert

- Methods based on obfuscating scene / image geometry
- Original geometry can be recovered quite easily if neighborhood information is available
- Neighborhoods can be approximately recovered from descriptors

See [Chelani, Benbihi et al., arXiv:2409.11536] for details

al., CVPR 2019]
CVPR 2023]
ICCV 2023]

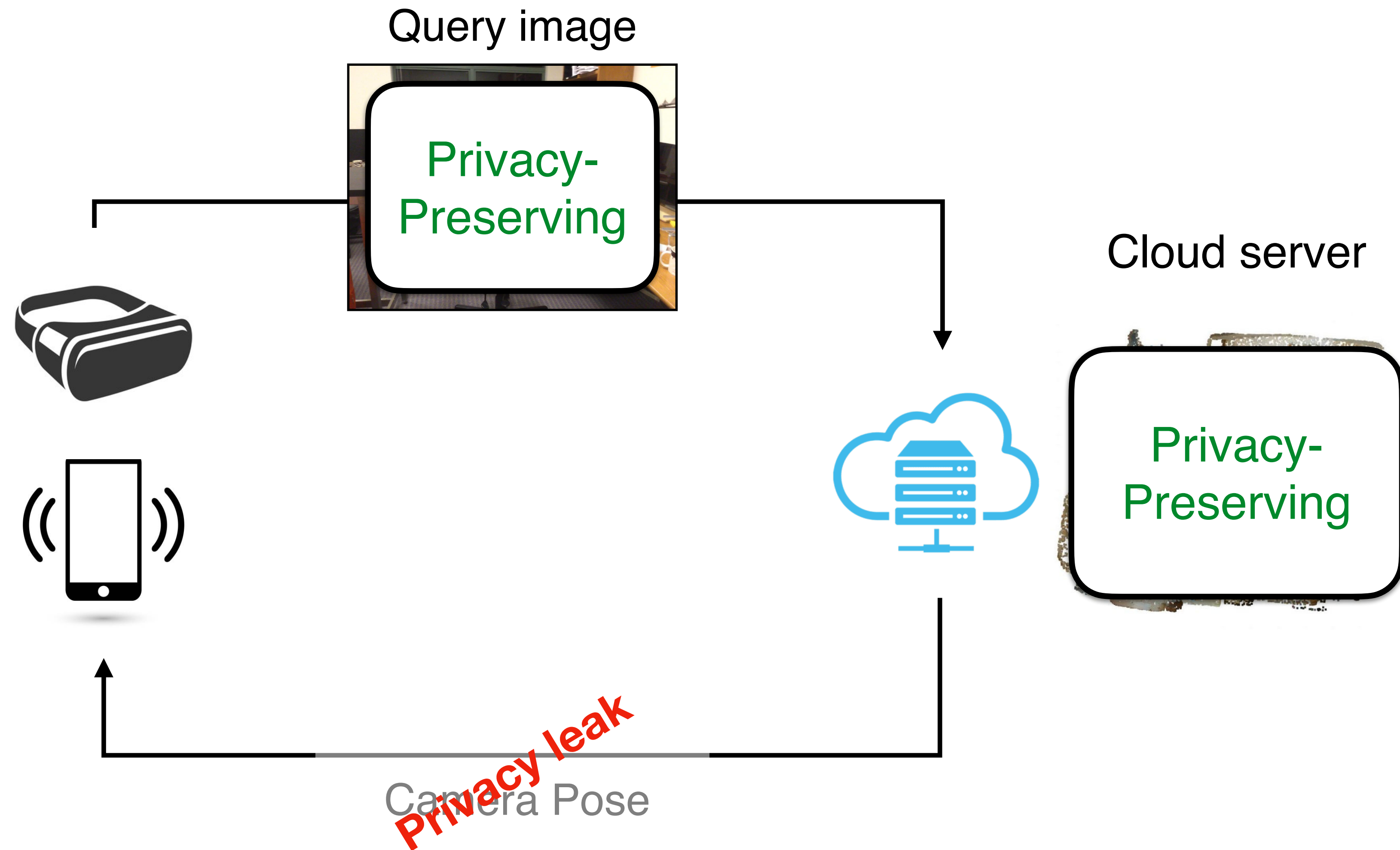
vate
y users

Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

[Chelani, Benbihi, Kahl, Sattler, Kukulova, Obfuscation Based Privacy Preserving Representations are Recoverable Using Neighborhood Information, arXiv:2409.11536]

Are Privacy-Preserving Representations Enough?



slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

The Downside of Robust Localization



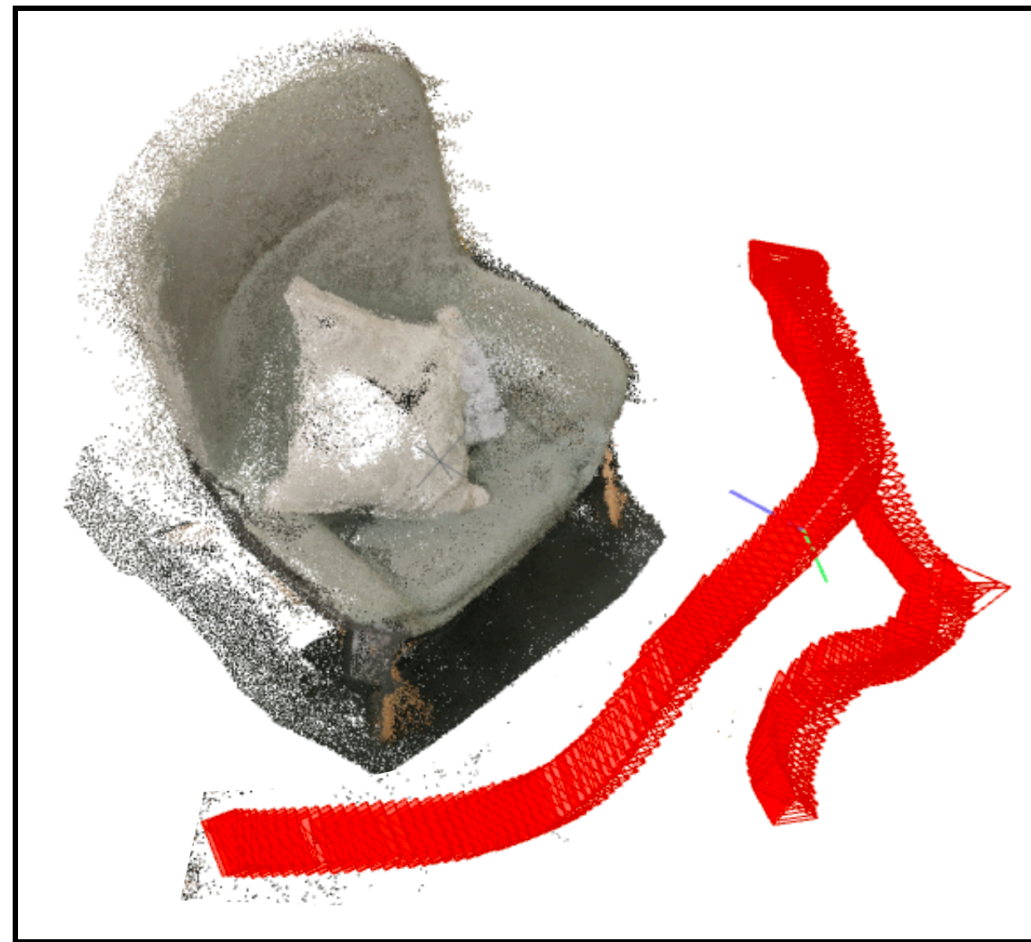
Robustness to shape and appearance variations means we can match images of different object instances

slide credit: Mihai Dusmanu, Kunal Chelani

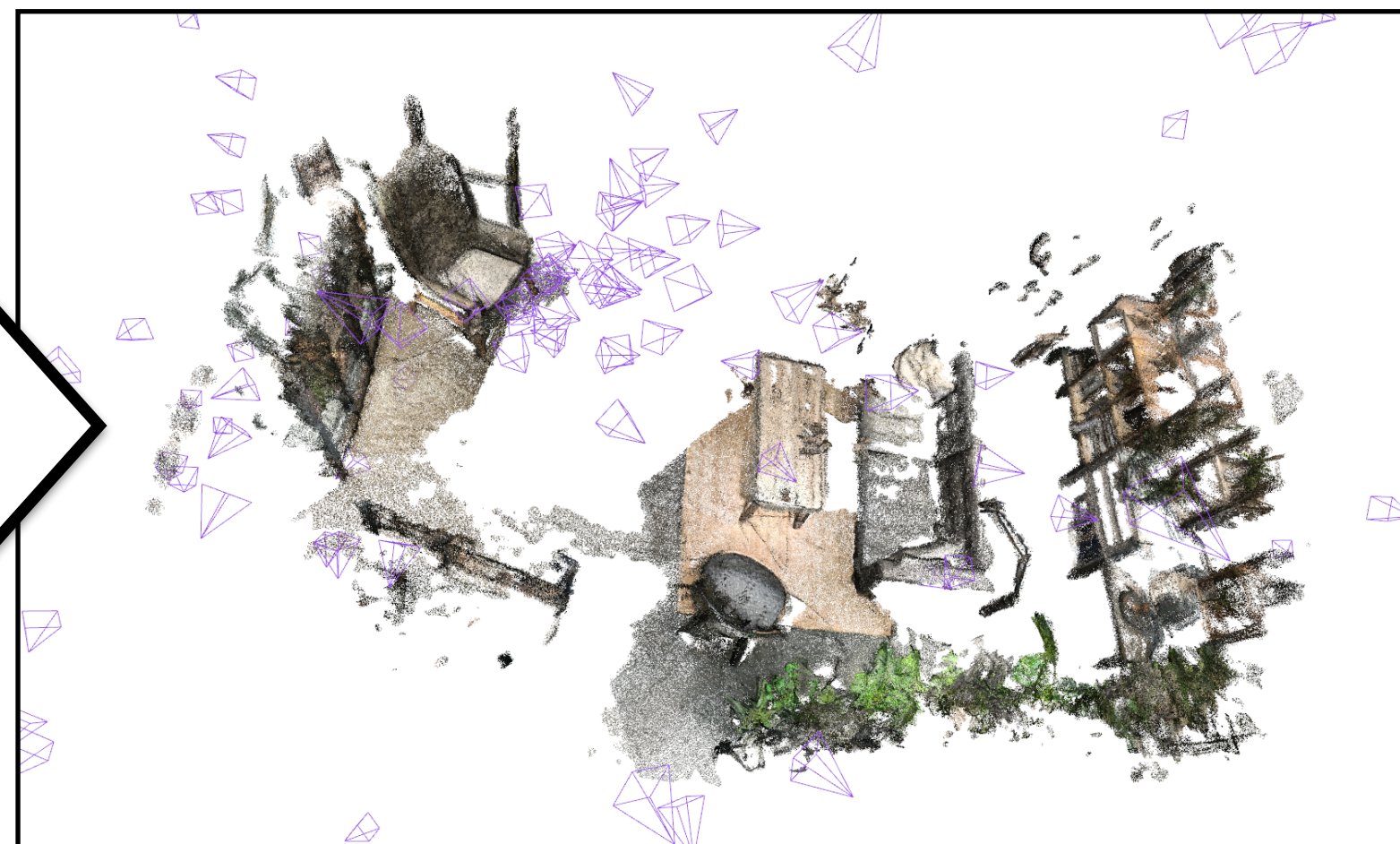
[Dusmanu, Rocco, Pajdla, Pollefeys, Sivic, Torii, Sattler, D2-Net: A Trainable CNN for Joint Detection and Description of Local Features, CVPR 2019]

Recovering Scene Content from Camera Poses

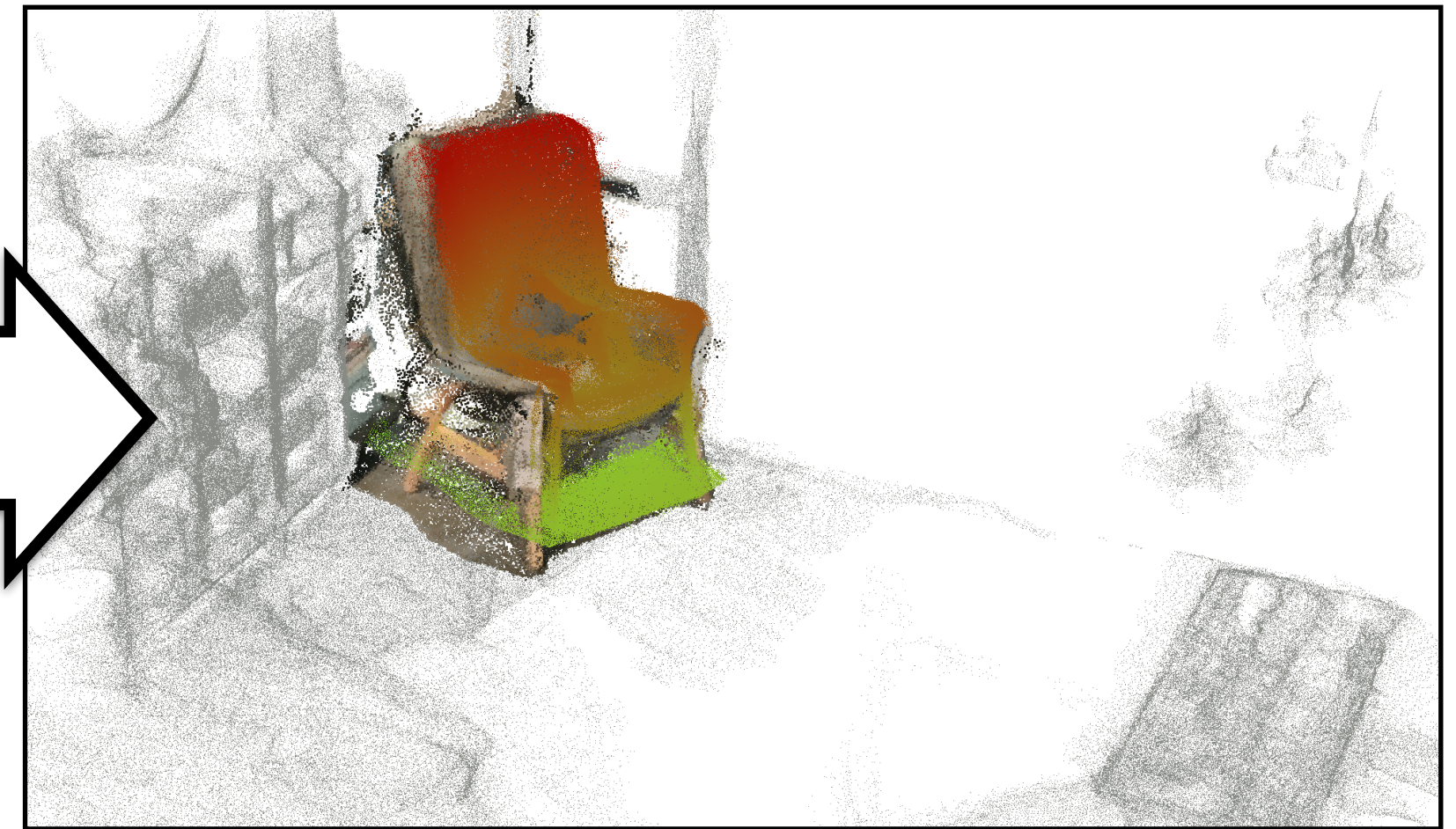
input:
image sequence of object
(e.g., from Internet)



Attacker runs SfM to
get camera poses



Camera poses returned by
the server



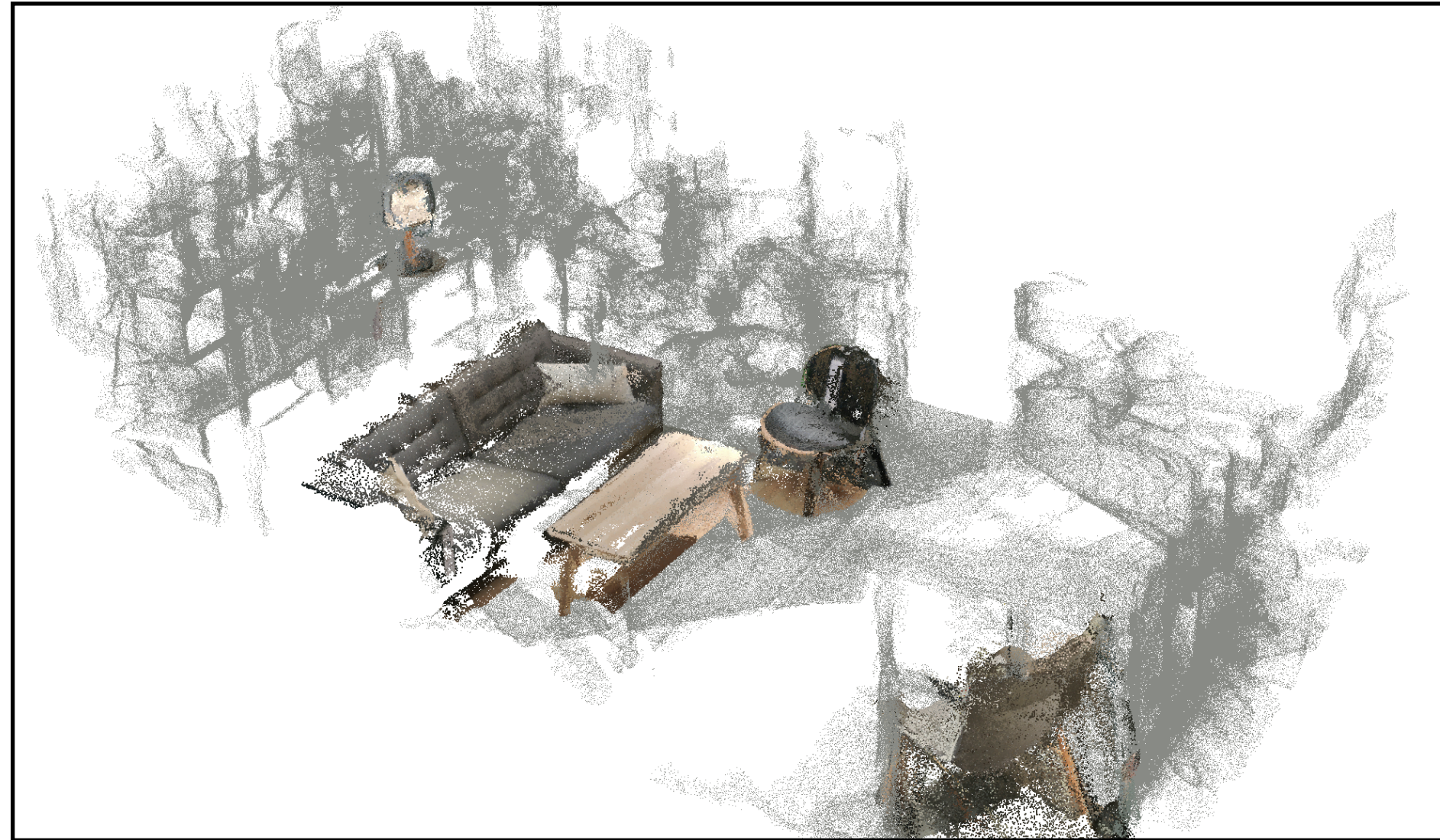
Object position from pose
alignment

slide credit: Kunal Chelani

[Chelani, Sattler, Kahl, Kukulova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

Qualitative Results

slide credit: Kunal Chelani



Actual scene with highlighted objects



Roughly reconstructed scene



[Chelani, Sattler, Kahl, Kukelova. Privacy-Preserving Representations are not Enough: Recovering Scene Content from Camera Poses, CVPR 2023]

Open Positions



- Open PhD & PostDoc position on camera geometry estimation, starting in 2025
- Contact: kukelova@gmail.com



- One open postdoc positions on privacy-preserving / temporal 3D mapping
- Open PhD position on visual localization
- Contact: torsten.sattler@cvut.cz