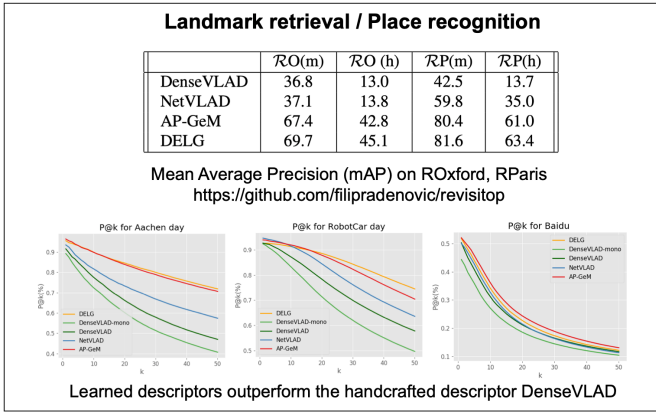
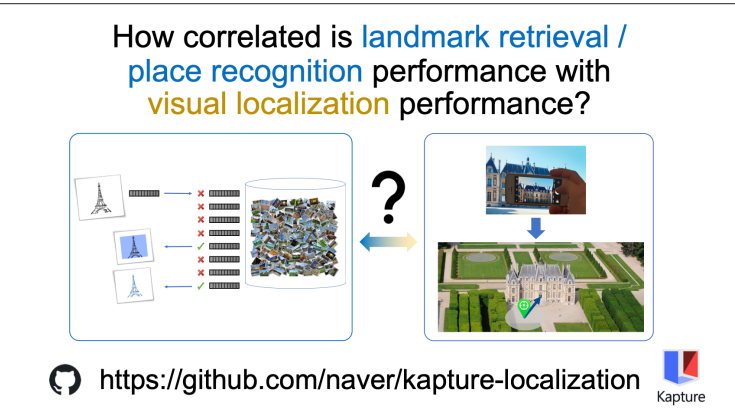
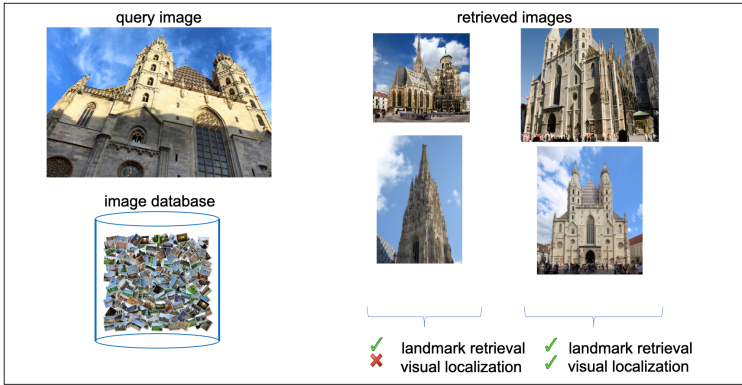


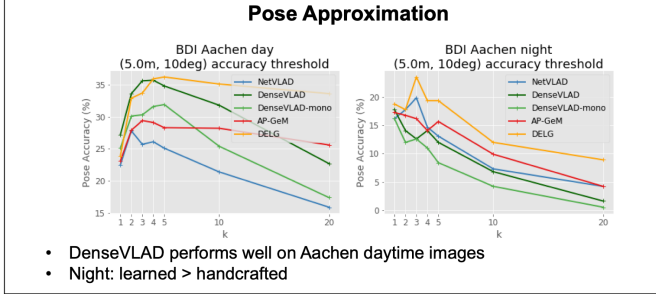
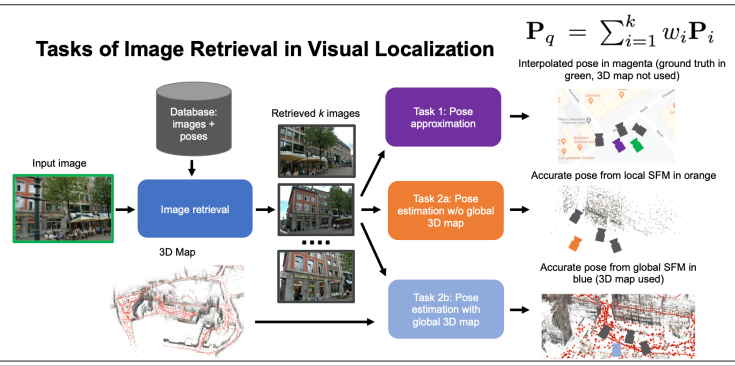
<sup>1</sup>Noé Pion, <sup>1</sup>Martin Humenberger, <sup>1</sup>Gabriela Csurka, <sup>1</sup>Yohann Cabon, <sup>2</sup>Torsten Sattler  
<sup>1</sup>NAVER LABS Europe, <sup>2</sup>Czech Technical University in Prague

This work received funding through the EU Horizon 2020 project RICAIP (grant agreement No 857306) and the European Regional Development Fund under IMPACT No. CZ.02.1.01/0.0/0.0/15 003/0000468.



### Benchmark

Aachen Day-Night-v1.1 [75]	Baidu-Mall [85]	RobotCar Seasons [56]
handheld scenario, augmented and mixed reality application	indoor	Autonomous driving
outdoor	close-up reflections and transparent surfaces, moving people	outdoor
day-night, various viewpoints	day-night, seasons, low image quality, dynamic traffic scenes	day-night, seasons, low image quality, dynamic traffic scenes



**Handcrafted global features:**

- DenseVLAD aggregates RootSIFT into an intra-normalized VLAD [90]

**Learned global features:**

- NetVLAD aggregates mid-level convolutional features extracted using VLAD [1]
- AP-GeM uses a generalized-mean pooling layer (GeM) to aggregate CNN-based descriptors [67]
- DELG is trained to extract both local and global features using one CNN [13]

**Handcrafted local features:**

- SIFT [53]

**Learned local features:**

- R2D2 [68]
- D2-Net [25]

**Evaluation:** <https://visuallocalization.net>

**Metrics:**

Visual localization accuracy:

- low (5m, 10°)
- medium (0.5m, 5°)
- high (0.25m, 2°)

Landmark retrieval / Place recognition

- P@k: 1 if all retrieved imgs are relevant
- R@k: 1 if among the k imgs at least one is relevant

[1] Arandjelovic et al. NetVLAD: CNN Architecture for Weakly Supervised Place Recognition, CVPR'16  
 [13] Cao et al. Unifying Deep Local and Global Features for Efficient Image Search, arXiv'20  
 [25] Dusmanu et al. D2-Net: A Trainable CNN for Joint Description and Detection of Local Features, CVPR'19  
 [53] Lowe, Distinctive Image Features from Scale-Invariant Keypoints, IJCV'04  
 [67] Revaud et al. Learning with Average Precision: Training Image Retrieval with a Listwise Loss, ICCV'19  
 [68] Revaud et al. R2D2: Reliable and Repeatable Detectors and Descriptors, NeurIPS, 2019  
 [56] Maddern et al. The Oxford RobotCar Dataset, IJRR'17  
 [75] Sattler, et al. Benchmarking 6DoF Outdoor Visual Localization in Changing Conditions, CVPR'18  
 [85] Sun, A Dataset for Benchmarking Image-based Localization, CVPR'17  
 [90] Torii et al. 24/7 Place Recognition by View Synthesis CVPR'15

