

Out with the Old?

CNN vs. SIFT-based visual localization

Prof. Dr. Laura Leal-Taixé

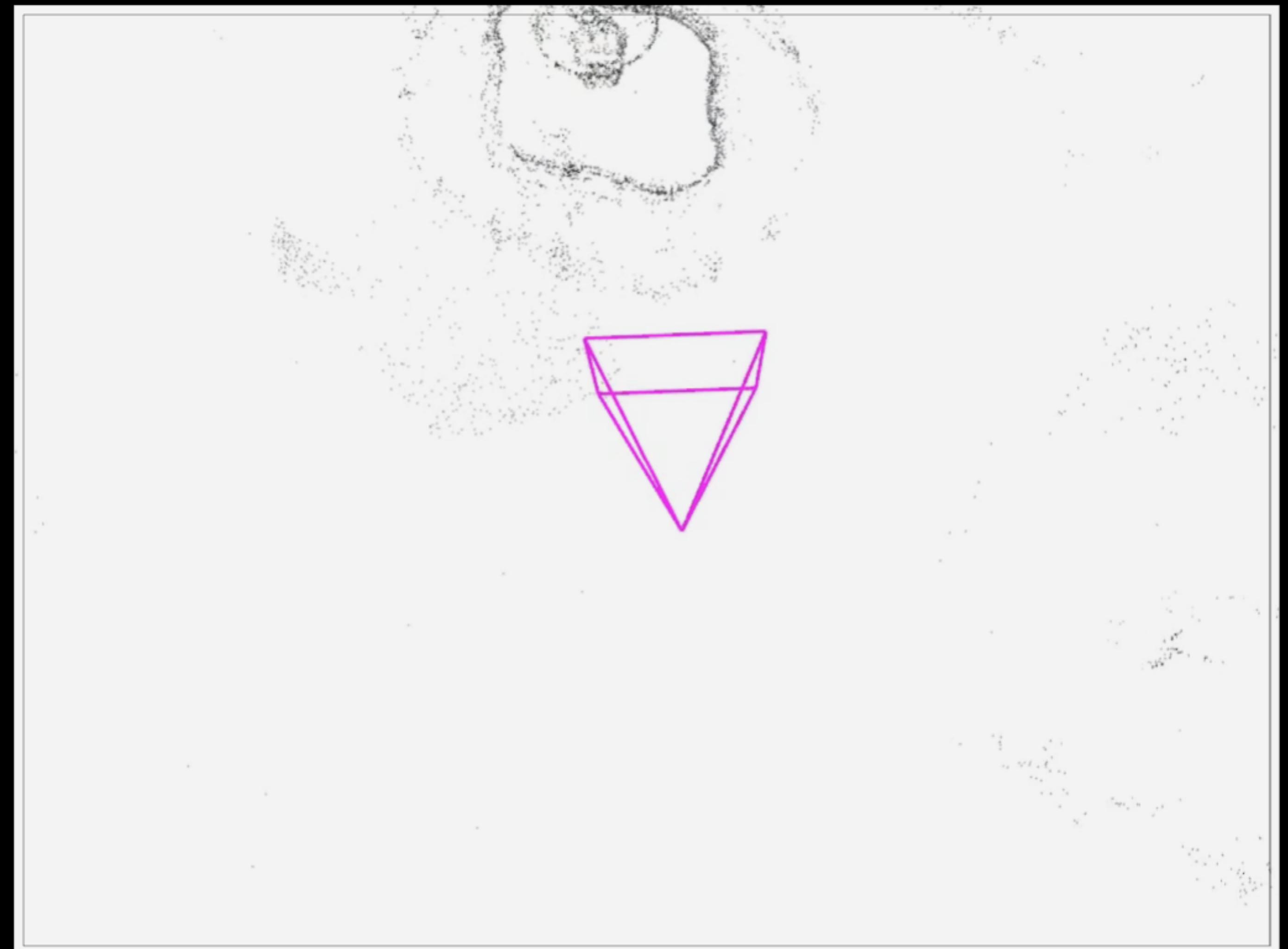


Technische Universität München

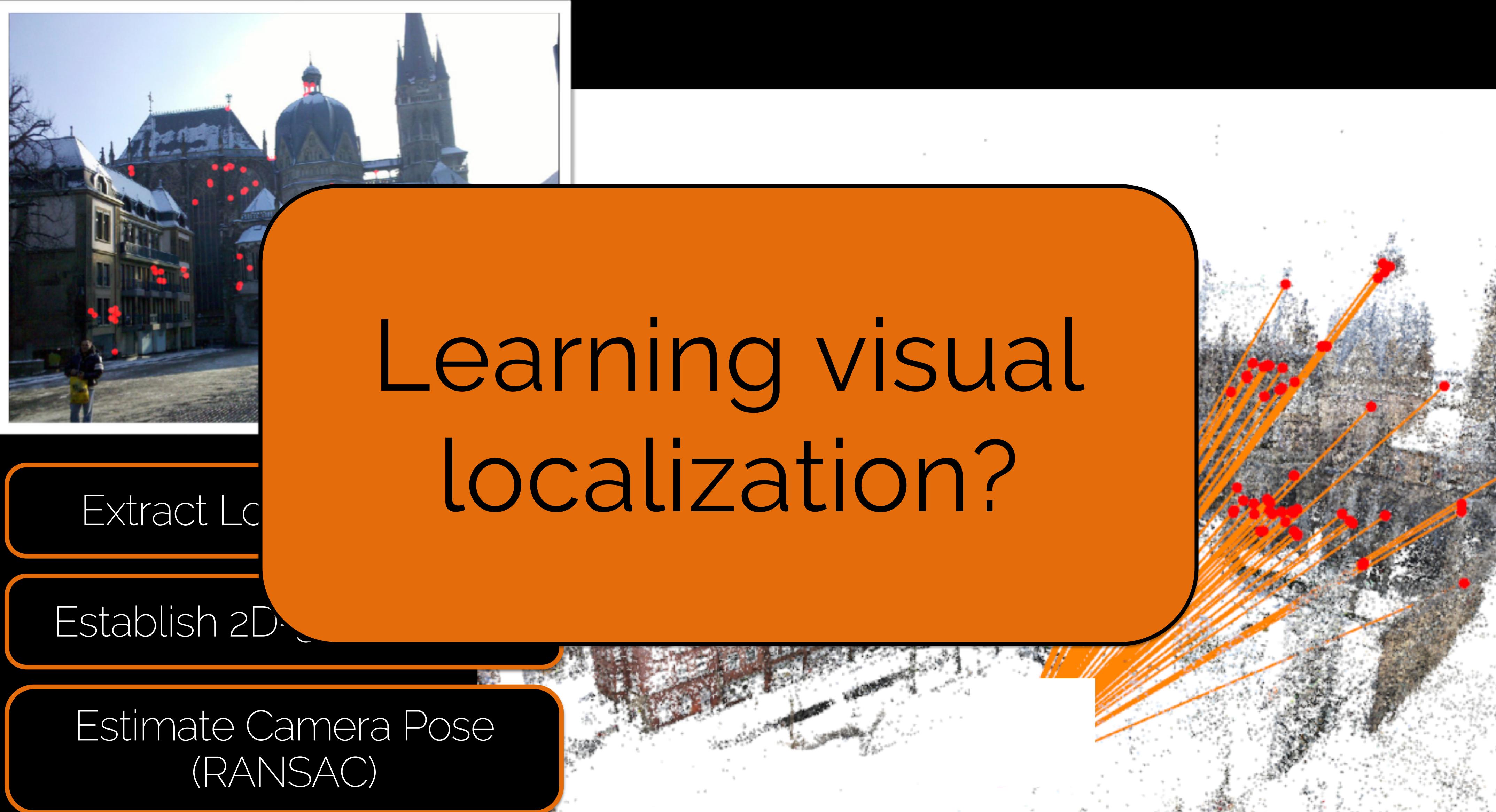
Visual Localization



Compute exact position and orientation of query image
(relative to 3D scene model)



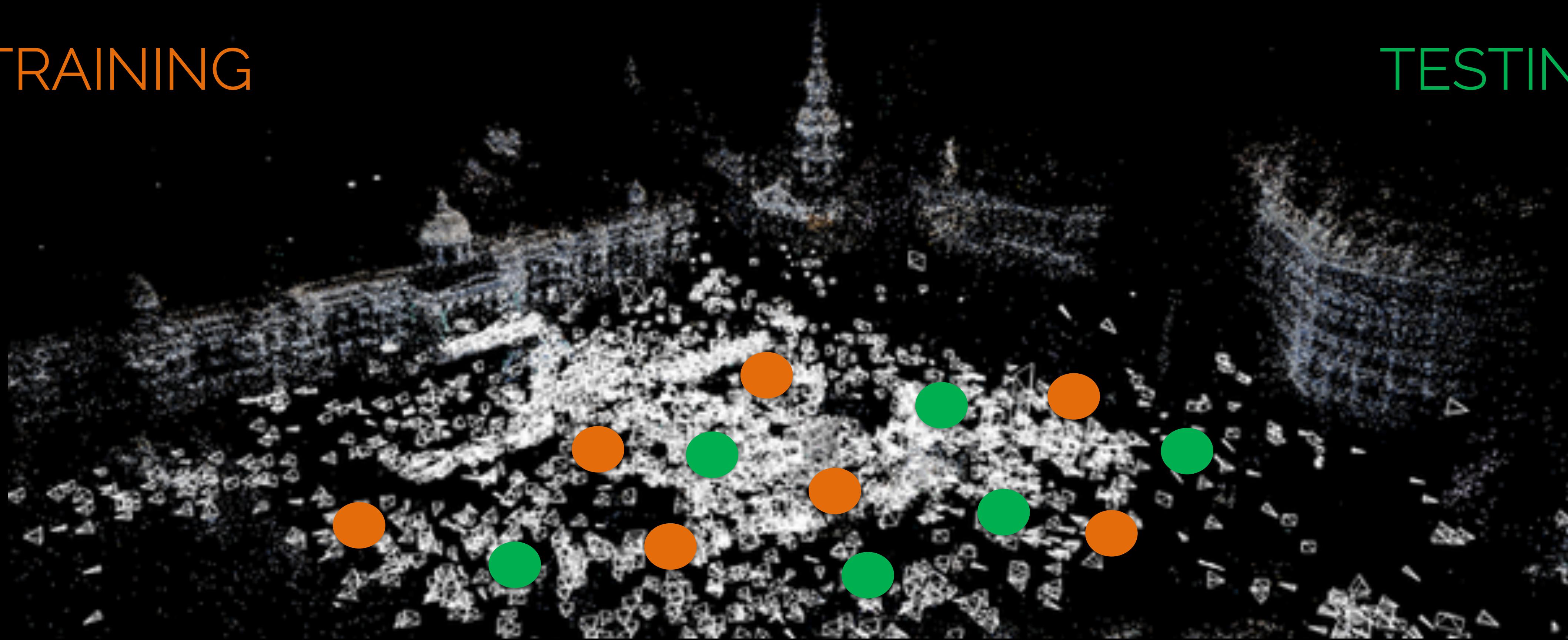
Classic Localization Pipeline



Where do we get training data?

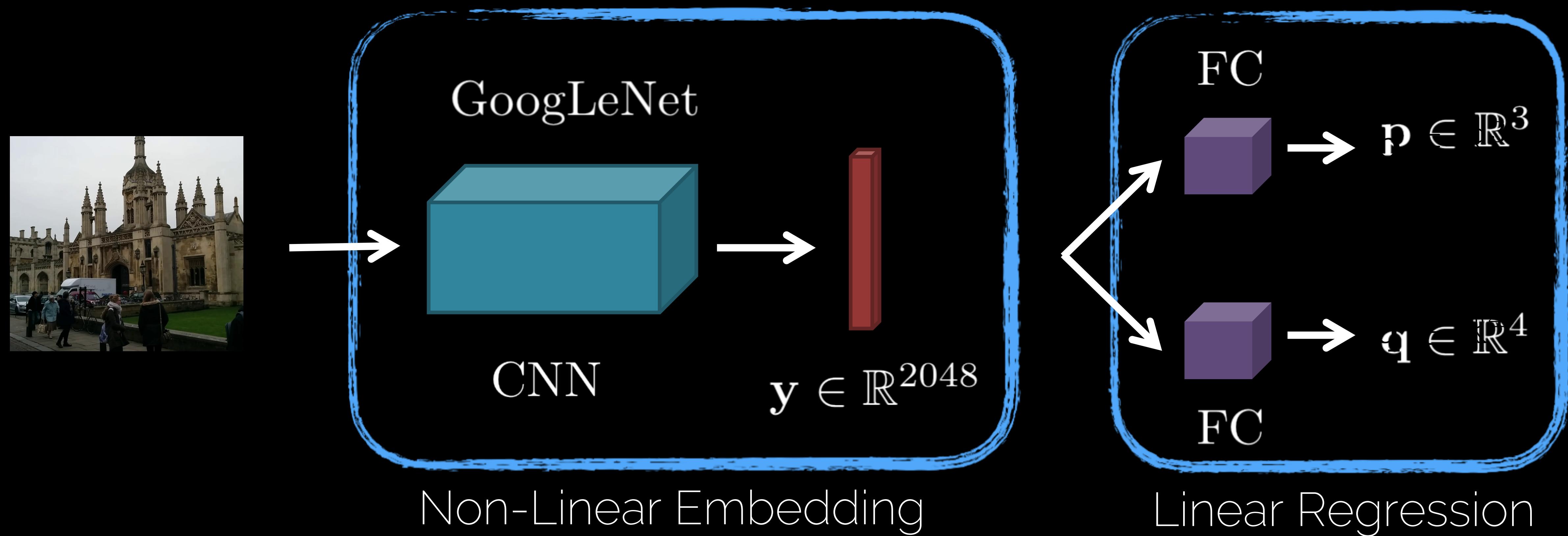
TRAINING

TESTING



Structure-from-Motion gives us plenty of data

Related work: PoseNet



Outdoor localization

- Experiments on Cambridge Landmarks

		CNN
Scene	Area or Volume	PoseNet [22]
King's College	5600 m^2	$1.92 \text{ m}, 5.40^\circ$
Old Hospital	2000 m^2	$2.31 \text{ m}, 5.38^\circ$
Shop Façade	875 m^2	$1.46 \text{ m}, 8.08^\circ$
St Mary's Church	4800 m^2	$2.65 \text{ m}, 8.48^\circ$
Average All	-	$2.08 \text{ m}, 6.83^\circ$
Average by [41]	-	-

Indoor localization: SIFT suffers

- Experiments on 7Scenes

Scene	Area or Volume	SIFT	CNN
		Active Search (w/o) [41]	PoseNet [22]
Chess	$6 m^3$	0.04 m, 1.96° (0)	0.32 m, 8.12°
Fire	$2.5 m^3$	0.03 m, 1.93° (1)	0.47 m, 14.4°
Heads	$1 m^3$	0.03 m, 2.59° (7)	0.29 m, 12.0°
Office	$7.5 m^3$	0.09 m, 3.61° (34)	0.48 m, 7.68°
Pumpkin	$5 m^3$	0.08 m, 3.10° (71)	0.47 m, 8.42°
Red Kitchen	$18 m^3$	0.07 m, 3.37° (0)	0.59 m, 8.64°
Stairs	$7.5 m^3$	not available	0.47 m, 13.8°
Average All	-	-	0.44 m, 10.4°
Average by [41]	-	0.06 m, 2.76°	-

Number of images that it failed to localize

Our new dataset TUM-LSI: SIFT dies

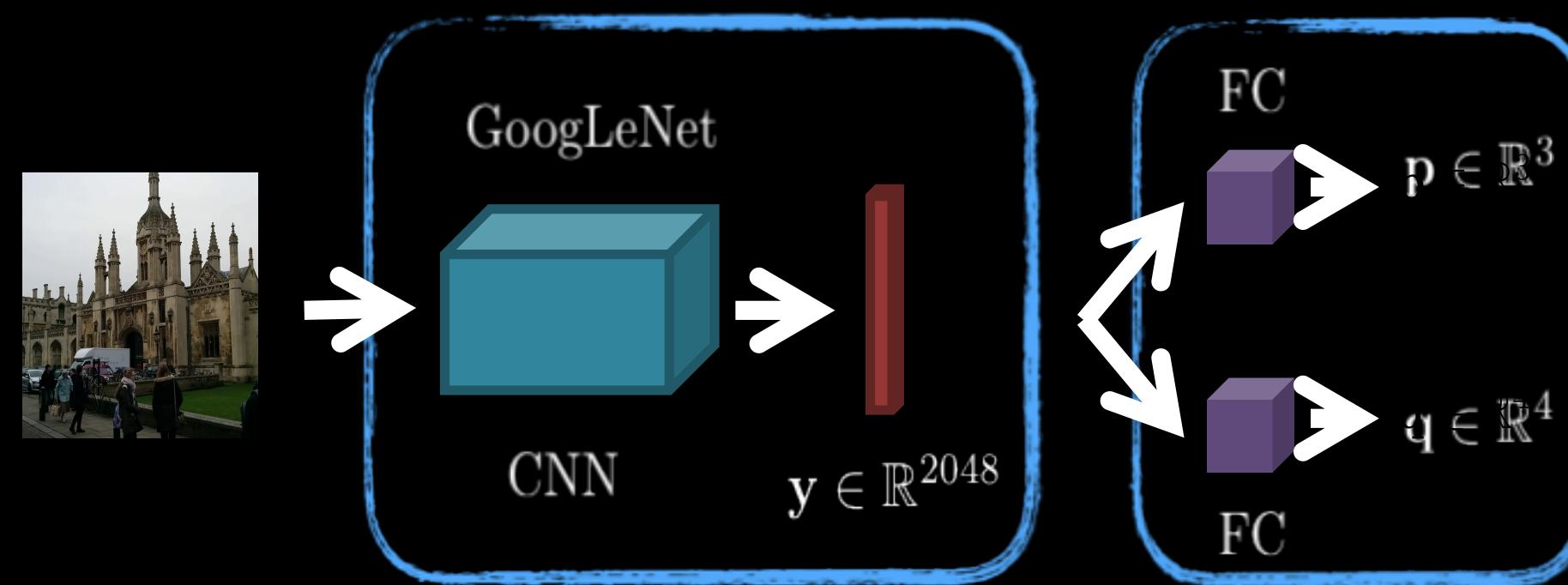
Scene	Area	# train/test	PoseNet [22]
Matriculation Hall	793 m ²	940/235	1.15 m, 3.93°
Hallway	2677 m ²	845/215	2.16 m, 6.99°
Floor	5575 m ²	875/220	1.87 m, 6.14°
Average			1.72 m, 5.68°



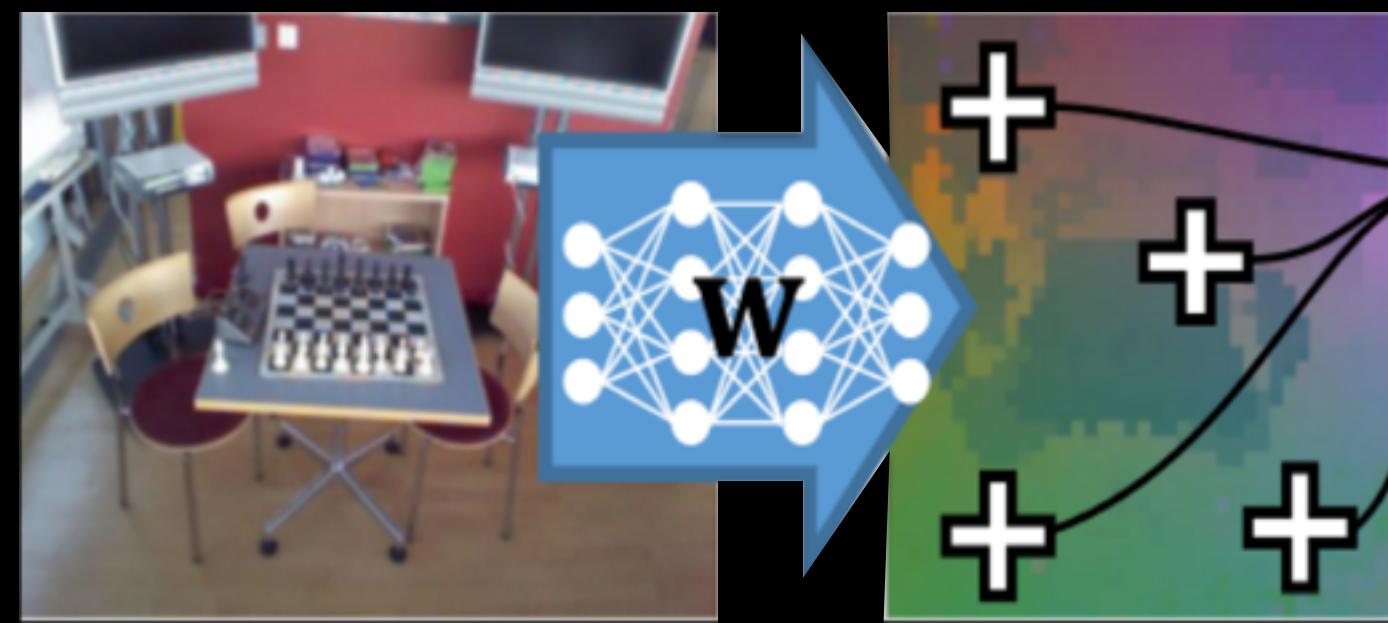
SIFT-based methods do not work at all!

Limitations of current methods

- One separate model is needed per scene



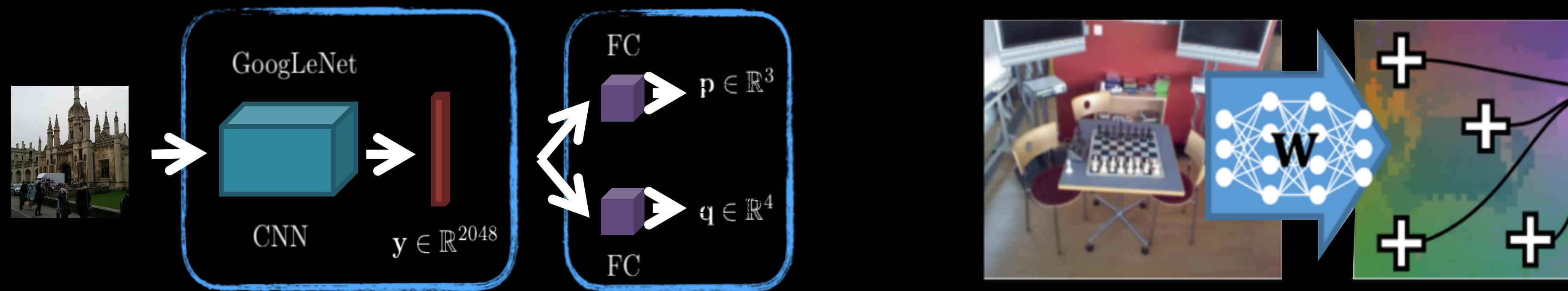
Kendall et al. PoseNet;. ICCV 2015



Brachmann et al. DSAC;. CVPR 2017

Limitations of current methods

- One separate model is needed per scene



- The most used loss function requires the setting of a hyperparameter which is again scene dependent

DOES NOT SCALE

PoseNet loss functions

- L2 loss

$$\mathcal{L}(y^*, y) = \| c - c^* \|_2 + \beta \| q - q^* \|_2$$

pose label

pose prediction

camera center
position

camera rotation
in quaternion

✗ No geometry information

scaling factor to balance
the weighting between
the rotation error and the
position error

✗ Need to find its value per scene → can range from 200 to 2500

PoseNet loss functions

- Geometric loss function: Minimize re-projection error of 3D points visible in image

- ✓ More accurate results
- ✗ Network needs to be pre-trained on L2 loss
- ✗ One network per scene

Relative Pose Estimation

- Use a neural network to predict relative poses
- Loss function: PoseNet L2 loss

$$\mathcal{L}(y^*, y) = \| c - c^* \|_2 + \boxed{\beta} \| q - q^* \|_2$$

- ✓ One network for all scenes
- ✗ Loss still depends on a scene-dependent hyperparameter

Limitations of current methods

- One separate model is needed per scene
- The most used loss function requires the setting of a hyperparameter which is again scene dependent

Proposed method

- One separate model is needed per scene

Relative Pose Estimation
with
Geometric Matching

- The most used loss function requires the setting of a hyperparameter which is again scene dependent

Essential Matrix
Regression

✓ One network
for all scenes

✓ No
hyperparameter
on the loss

Proposed method

- One separate model is needed per scene

Relative Pose Estimation
with
Geometric Matching

- The most used loss function requires the setting of a hyperparameter which is again scene dependent

Essential Matrix
Regression

✓ One network
for all scenes

✓ No
hyperparameter
on the loss

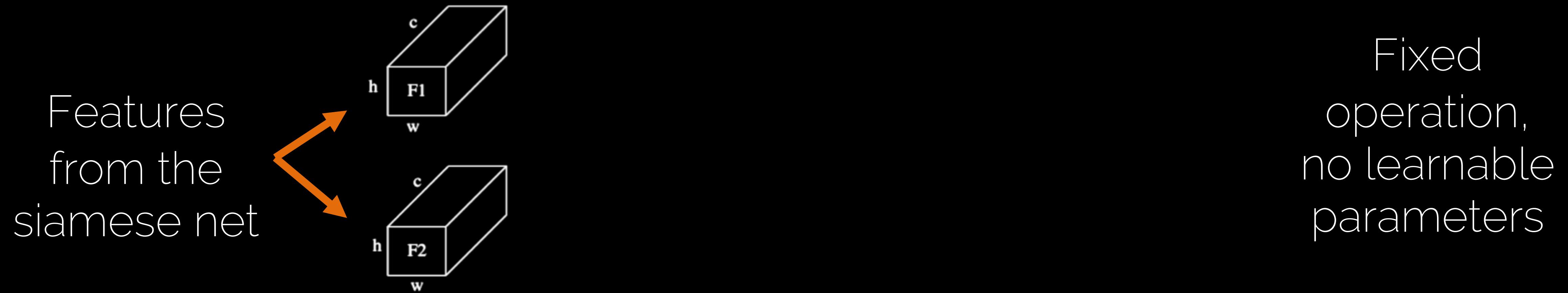
Regressing relative poses

- Siamese architecture
 - Shared weights on the input image pair
 - Feature concatenation



Geometric matching layer

- Mimic the geometrical matching process



Proposed method

- One separate model is needed per scene

Relative Pose Estimation
with
Geometric Matching

- The most used loss function requires the setting of a hyperparameter which is again scene dependent

Essential Matrix
Regression

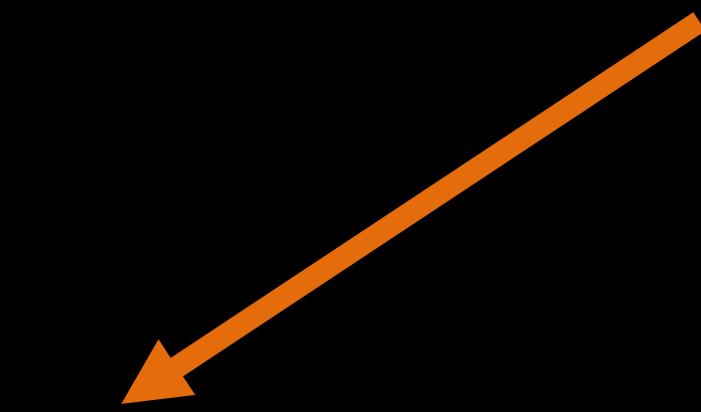
✓ One network
for all scenes

✓ No
hyperparameter
on the loss

Essential Matrix loss

- Essential matrix between query and training image $\mathbf{E}^* = [\mathbf{t}^*]_\times \mathbf{R}^*$
- L2 loss function on the entries of the essential matrix

$$\mathcal{L}_{ess}(\mathbf{E}^*, \mathbf{E}) = \| \mathbf{e} - \mathbf{e}^* \|_2$$

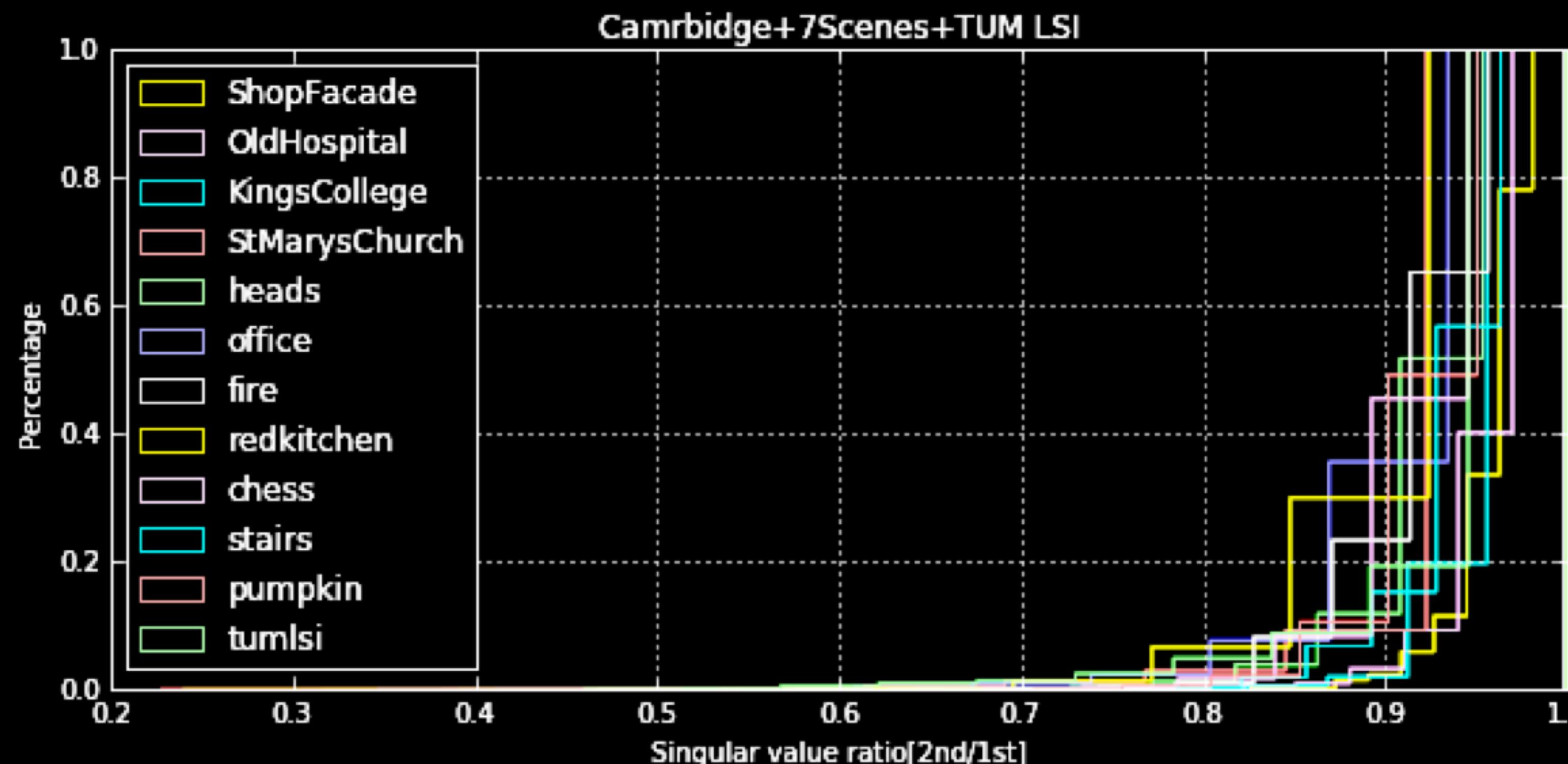
$$\| \mathbf{t} \| = 1$$


The loss is
hyperparameter
free!

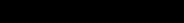
Fixing the scale of translations to
get consistent predictions

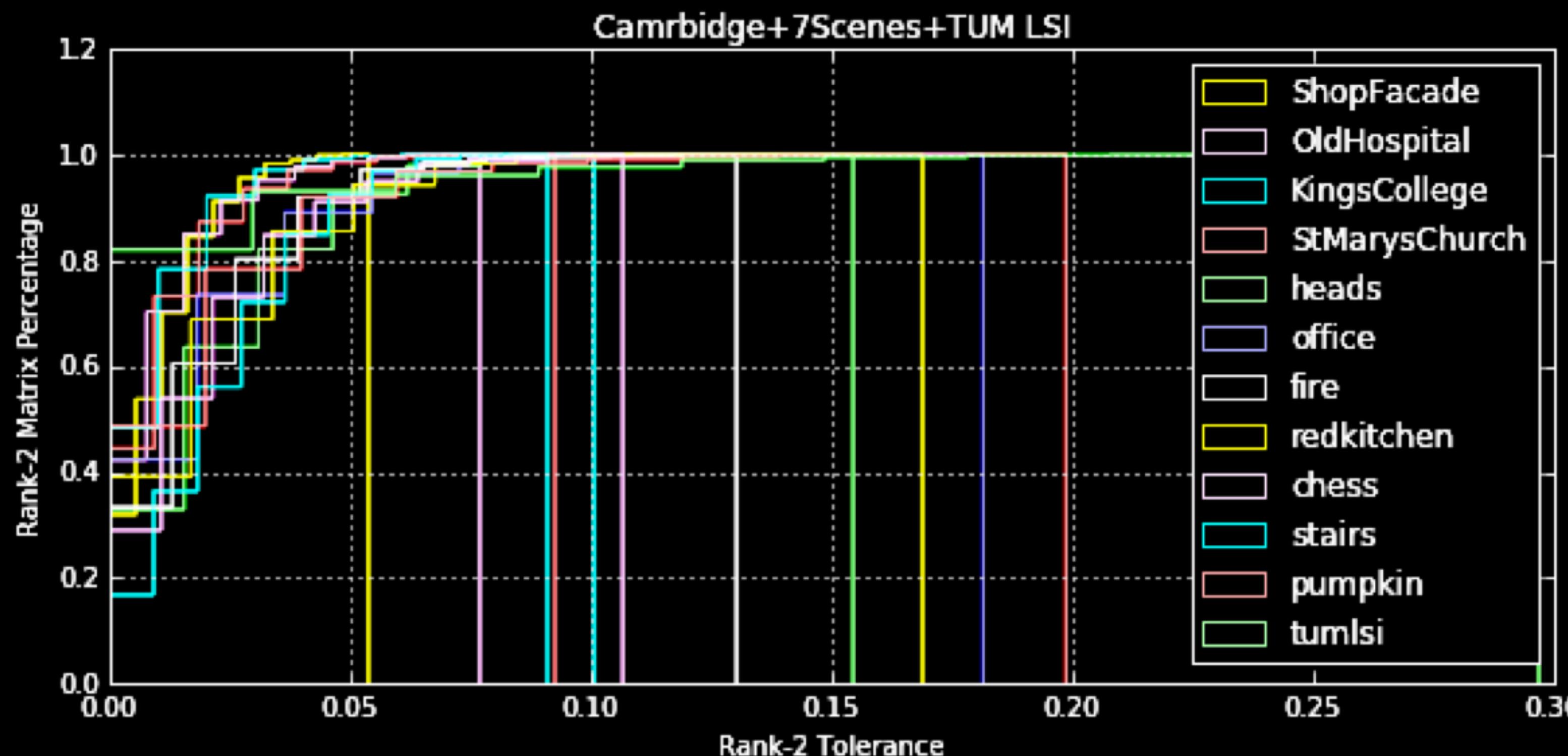
Are we really learning an essential matrix?

- Ratio between first two eigenvalues is close to 1? 



Are we really learning an essential matrix?

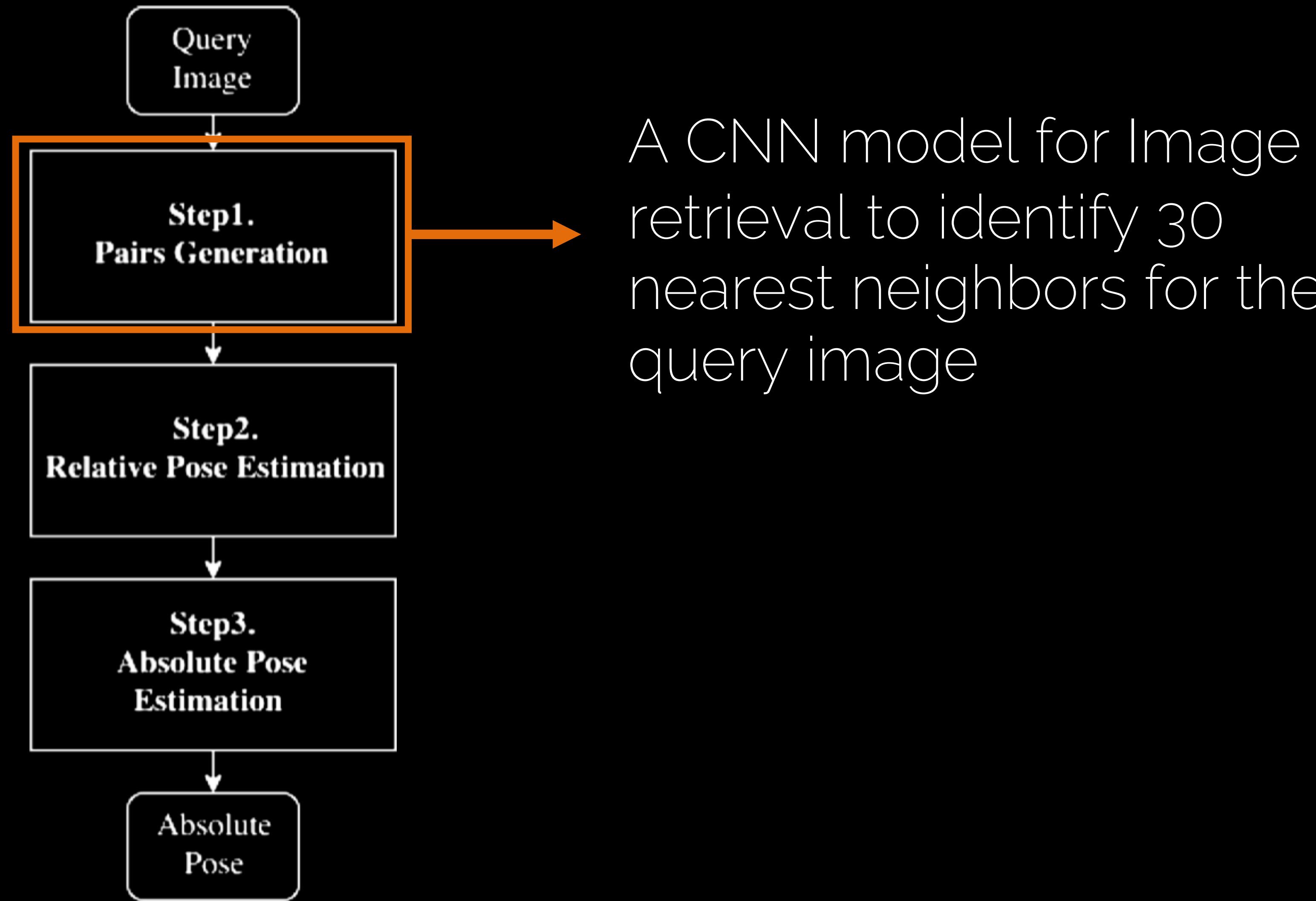
- Ratio between first two eigenvalues is close to 1? 
 - Third eigenvalue close to 0? 



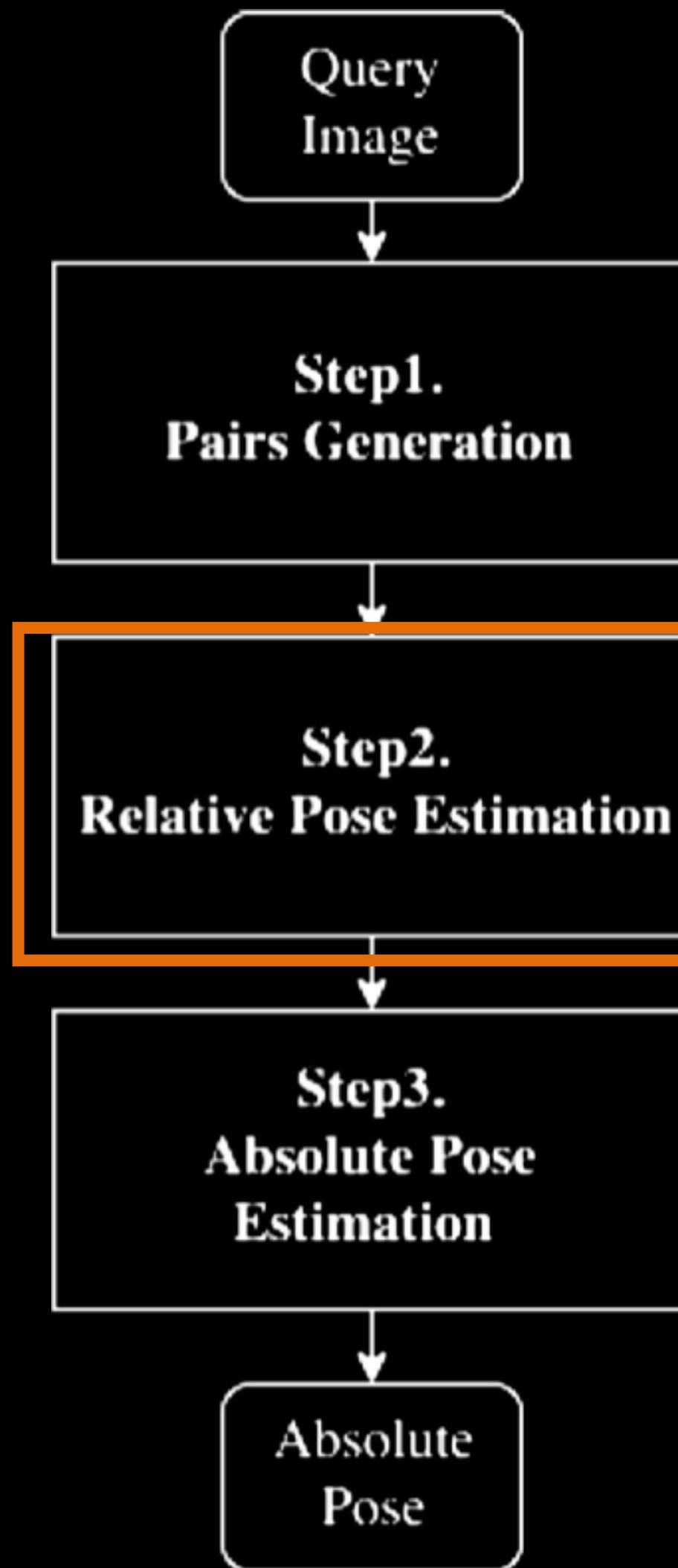
Are we really learning an essential matrix?

- Ratio between first two eigenvalues is close to 1? 
- Third eigenvalue close to 0? 
- We learn a good approximation of a true essential matrix
- Feature-based methods also project the resulting matrix to a rank 2 matrix (due to noise)
- Forcing this projection **does not have any effect** on the results on any of the tested datasets

Full localization pipeline



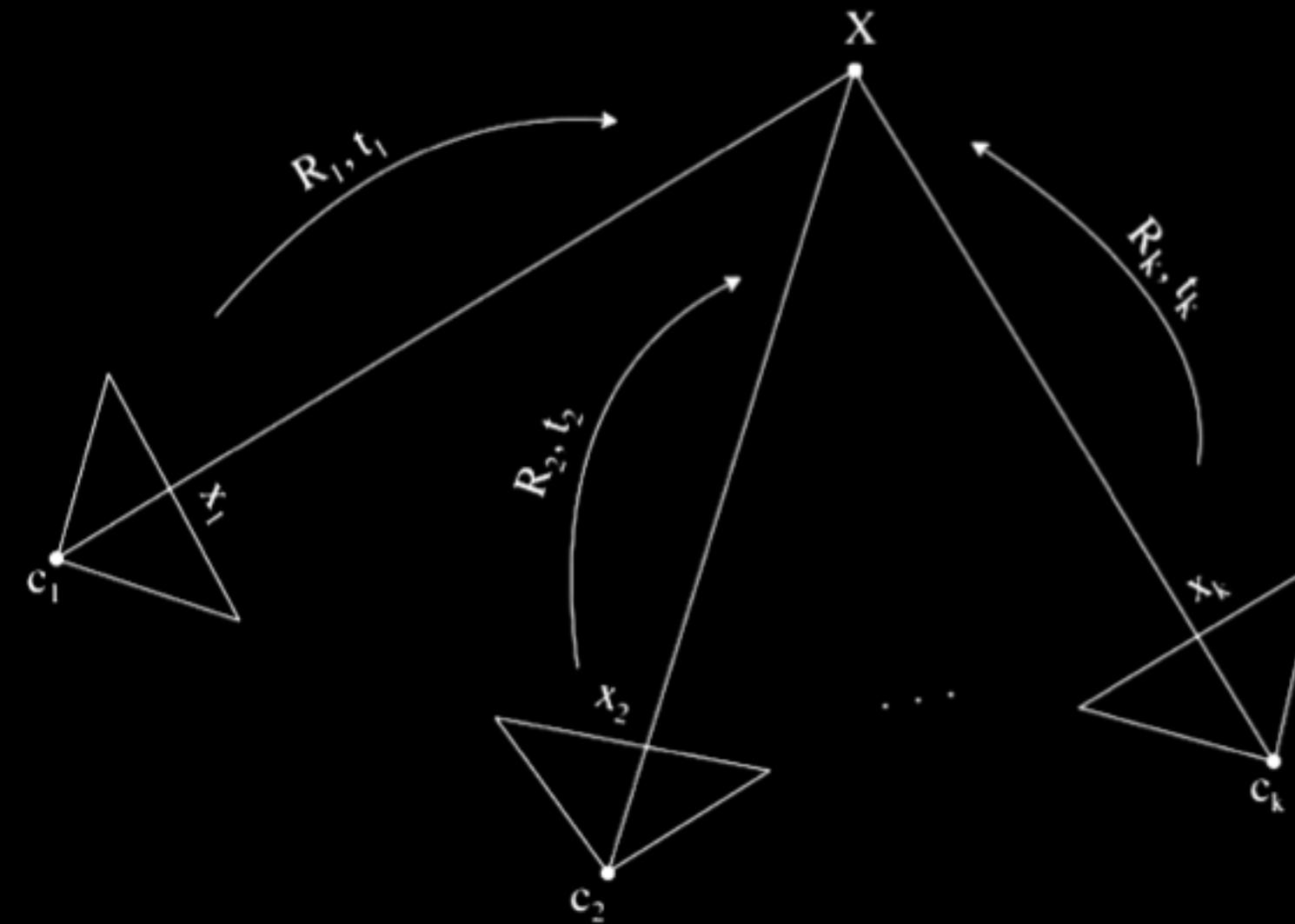
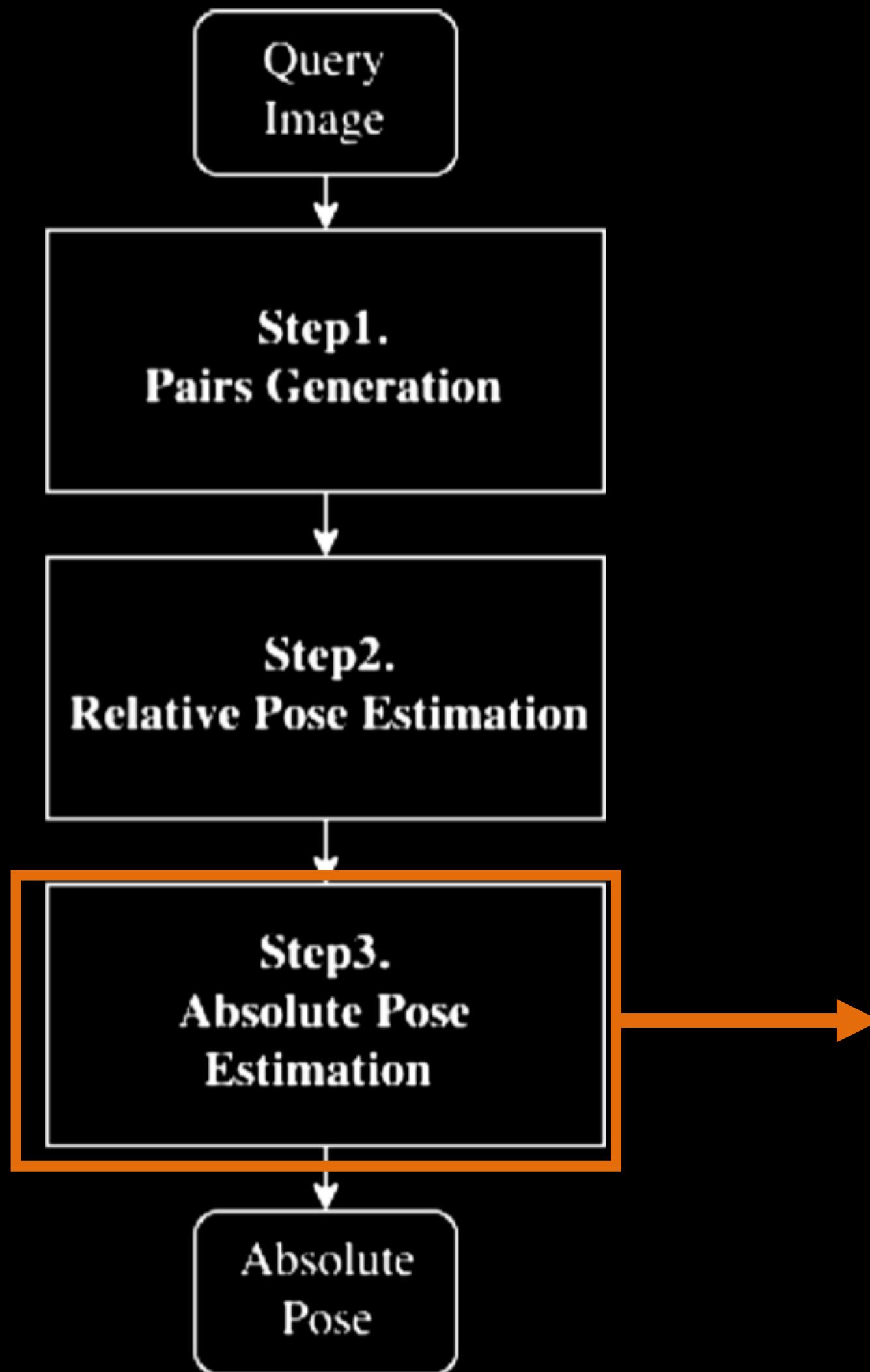
Full localization pipeline



Step2.
Relative Pose Estimation

EssNet to regress essential matrices and estimate relative poses

Full localization pipeline



Triangulation of the position and averaging of the rotation within a RANSAC loop.

Comparison to SOA

- Experiments on outdoor scenes, Cambridge Landmarks

Comparison to SOA

- Experiments on outdoor scenes, Cambridge Landmarks

Dataset	PoseNet [25]	Bayesian PoseNet [26]	LSTM PoseNet [28]	Geometric PoseNet [27]	VGGRegNet [32]	EssNet
St Mary's Church	$2.65m, 8.48^\circ$	$2.11m, 8.38^\circ$	$1.52m, 6.68^\circ$	$1.57m, 3.32^\circ$	$2.11m, 8.11^\circ$	$1.39m, 5.61^\circ$
ShopFacade	$1.46m, 8.08^\circ$	$1.25m, 7.54^\circ$	$1.18m, 7.44^\circ$	$0.88m, 3.78^\circ$	$0.63m, 5.73^\circ$	$0.87m, 4.58^\circ$
OldHospital	$2.31m, 5.38^\circ$	$2.57m, 5.14^\circ$	$1.51m, 4.29^\circ$	$3.20m, 3.29^\circ$	$1.50m, 4.03^\circ$	$2.14m, 5.84^\circ$
King's College	$1.92m, 5.40^\circ$	$1.74m, 4.06^\circ$	$0.99m, 3.65^\circ$	$0.88m, 1.04^\circ$	$1.06m, 2.81^\circ$	$0.97m, 2.80^\circ$
Average	$2.09m, 6.84^\circ$	$1.92m, 6.28^\circ$	$1.30m, 5.52^\circ$	$1.63m, 2.86^\circ$	$1.33m, 5.17^\circ$	$1.34m, 4.71^\circ$

- SOA positional error with better rotational error

Comparison to SOA

- Experiments on outdoor scenes, Cambridge Landmarks

Dataset	PoseNet [25]	Bayesian PoseNet [26]	LSTM PoseNet [28]	Geometric PoseNet [27]	VGGRegNet [32]	EssNet
St Mary's Church	$2.65m, 8.48^\circ$	$2.11m, 8.38^\circ$	$1.52m, 6.68^\circ$	$1.57m, 3.32^\circ$	$2.11m, 8.11^\circ$	$1.39m, 5.61^\circ$
ShopFacade	$1.46m, 8.08^\circ$	$1.25m, 7.54^\circ$	$1.18m, 7.44^\circ$	$0.88m, 3.78^\circ$	$0.63m, 5.73^\circ$	$0.87m, 4.58^\circ$
OldHospital	$2.31m, 5.38^\circ$	$2.57m, 5.14^\circ$	$1.51m, 4.29^\circ$	$3.20m, 3.29^\circ$	$1.50m, 4.03^\circ$	$2.14m, 5.84^\circ$
King's College	$1.92m, 5.40^\circ$	$1.74m, 4.06^\circ$	$0.99m, 3.65^\circ$	$0.88m, 1.04^\circ$	$1.06m, 2.81^\circ$	$0.97m, 2.80^\circ$
Average	$2.09m, 6.84^\circ$	$1.92m, 6.28^\circ$	$1.30m, 5.52^\circ$	$1.63m, 2.86^\circ$	$1.33m, 5.17^\circ$	$1.34m, 4.71^\circ$

- Only geometric loss is more accurate in rotation but 21% less accurate in translation

Comparison to SOA

- Experiments on outdoor scenes, Cambridge Landmarks

Dataset	PoseNet [25]	Bayesian PoseNet [26]	LSTM PoseNet [28]	Geometric PoseNet [27]	VGGRegNet [32]	EssNet
St Mary's Church	$2.65m, 8.48^\circ$	$2.11m, 8.38^\circ$	$1.52m, 6.68^\circ$	$1.57m, 3.32^\circ$	$2.11m, 8.11^\circ$	$1.39m, 5.61^\circ$
ShopFacade	$1.46m, 8.08^\circ$	$1.25m, 7.54^\circ$	$1.18m, 7.44^\circ$	$0.88m, 3.78^\circ$	$0.63m, 5.73^\circ$	$0.87m, 4.58^\circ$
OldHospital	$2.31m, 5.38^\circ$	$2.57m, 5.14^\circ$	$1.51m, 4.29^\circ$	$3.20m, 3.29^\circ$	$1.50m, 4.03^\circ$	$2.14m, 5.84^\circ$
King's College	$1.92m, 5.40^\circ$	$1.74m, 4.06^\circ$	$0.99m, 3.65^\circ$	$0.88m, 1.04^\circ$	$1.06m, 2.81^\circ$	$0.97m, 2.80^\circ$
Average	$2.09m, 6.84^\circ$	$1.92m, 6.28^\circ$	$1.30m, 5.52^\circ$	$1.63m, 2.86^\circ$	$1.33m, 5.17^\circ$	$1.34m, 4.71^\circ$

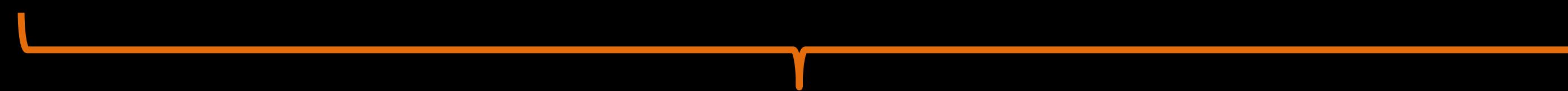
One network trained per scene

One single network

What can we do with more data?

- Training relative poses allows us to leverage more data

Dataset	1DSfM	Paris	1DSfM+Paris	
St Mary's Church	$1.75m, 6.88^\circ$	$1.26m, 6.28^\circ$	$1.19m, 5.13^\circ$	$1.39m, 5.61^\circ$
ShopFacade	$1.32m, 6.28^\circ$	$1.01m, 7.07^\circ$	$0.82m, 5.61^\circ$	$0.87m, 4.58^\circ$
OldHospital	$1.92m, 5.84^\circ$	$2.06m, 5.13^\circ$	$1.71m, 4.58^\circ$	$2.14m, 5.84^\circ$
King's College	$0.99m, 3.62^\circ$	$0.84m, 3.24^\circ$	$0.83m, 2.81^\circ$	$0.97m, 2.80^\circ$
Average	$1.50m, 5.66^\circ$	$1.29m, 5.43^\circ$	$1.14m, 4.53^\circ$	$1.34m, 4.71^\circ$



Pre-trained on a new dataset and fine-tuned on Cambridge Landmarks



Trained on Cambridge Landmarks

Out with the Old?

- SIFT-based methods still outperform CNN-based methods but there is potential in indoor scenes
- EssNet: relative pose estimation
 - mimics the visual localization pipeline → takes advantage of multiple view geometry
 - One network to perform localization in any scene
- Generalization still has a long way to go: if no images of the test scene are shown to the network, performance drops

Thank you!

Prof. Laura Leal-Taixé

<https://dvl.in.tum.de>



Technische Universität München