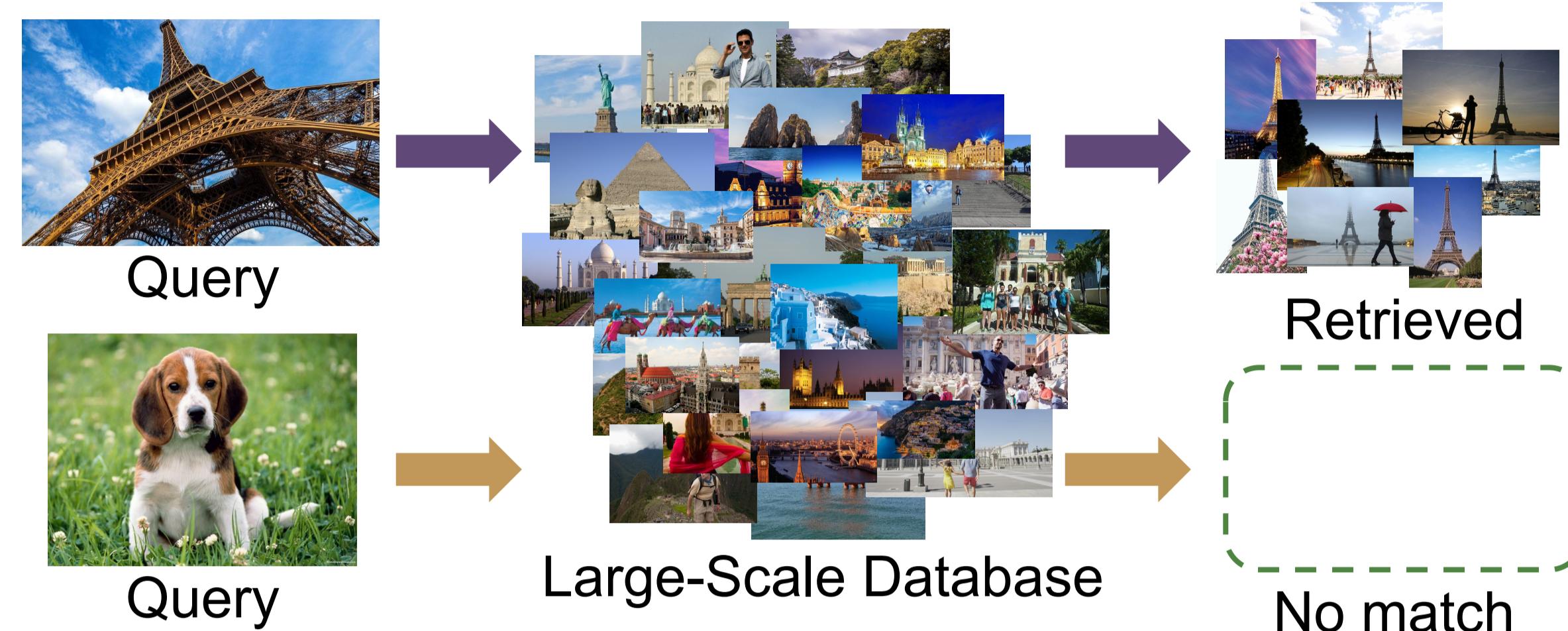


## Large-Scale Image Retrieval with Attentive Deep Local Features

Hyeonwoo Noh<sup>1</sup>, André Araujo<sup>2</sup>, Jack Sim<sup>2</sup>, Tobias Weyand<sup>2</sup>, Bohyung Han<sup>1</sup>[github.com/tensorflow/models/tree/master/research/delf](https://github.com/tensorflow/models/tree/master/research/delf) <sup>1</sup>POSTECH, Korea <sup>2</sup>Google Inc.

## Large-Scale Image Retrieval



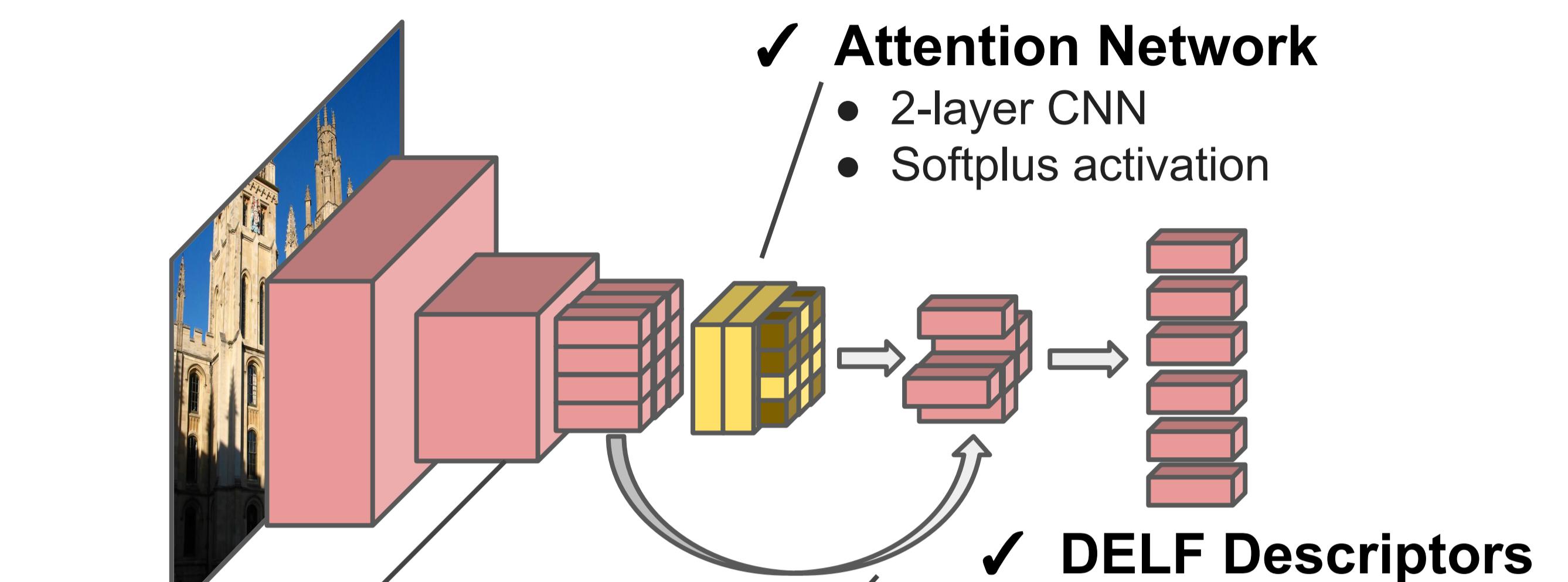
## Challenges

- Image clutter
- Partial occlusion
- Multiple landmarks
- Queries with no match
- Local features lack semantic information
- Patch-level annotations are expensive
- Existing dataset are small/medium

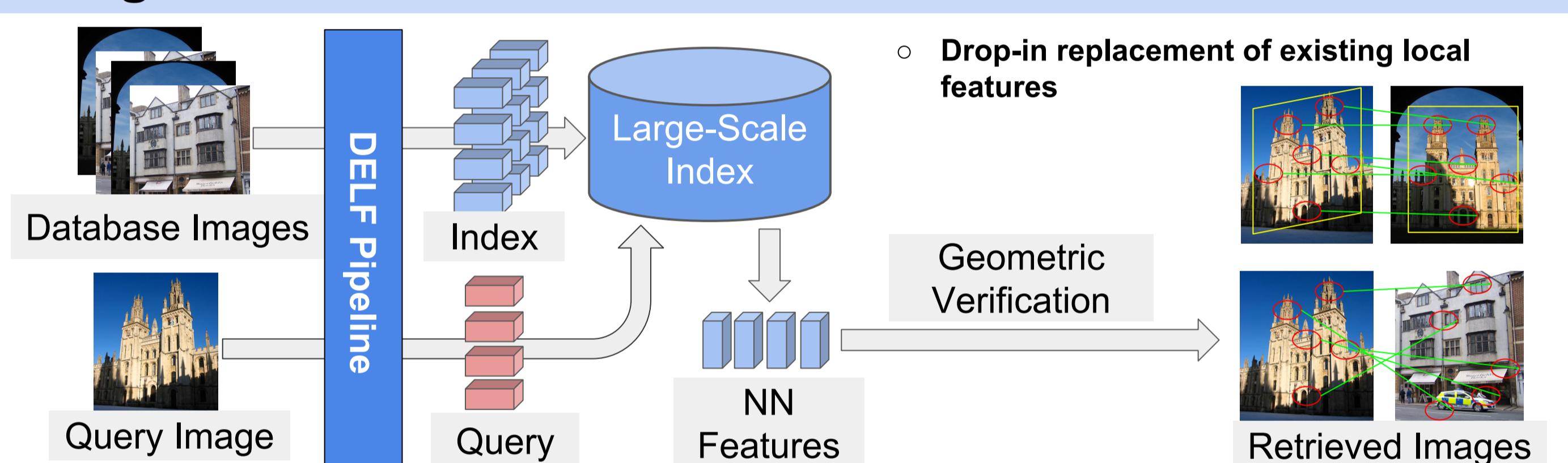
## Contributions

- Deep local feature retrieval with geometric verification
- High-level attention based keypoint selection
- Weakly-supervised feature learning
- New large-scale dataset

## DELF Pipeline



## Image Retrieval with DELF

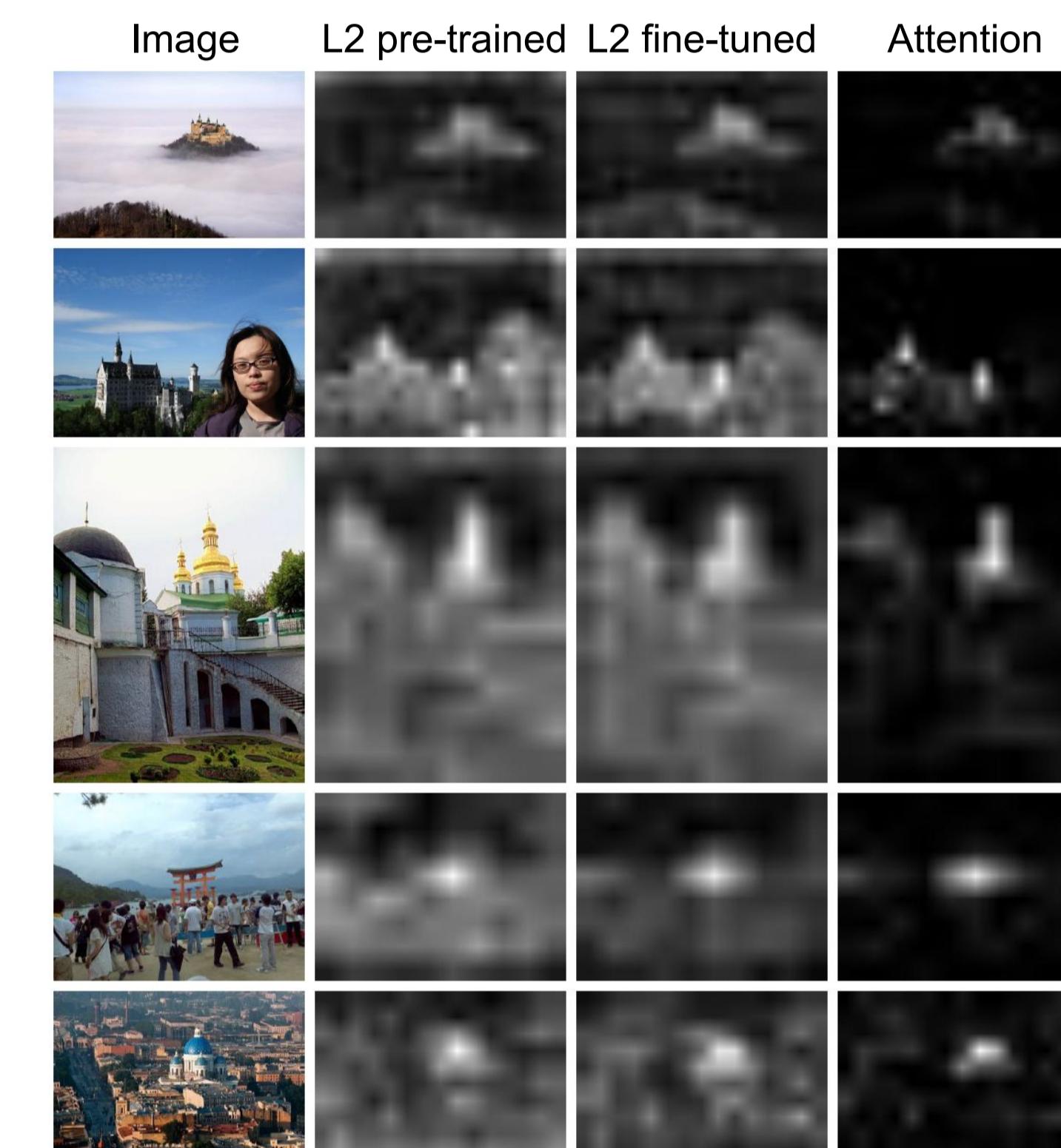


## DELF: DEep Local Features

## □ Attention-based Keypoint Selection

- Semantic keypoint selection based on high-level features
- Focus on discriminative features

## □ Attention Visualization



## □ DELF Learning

## Two-stage training:

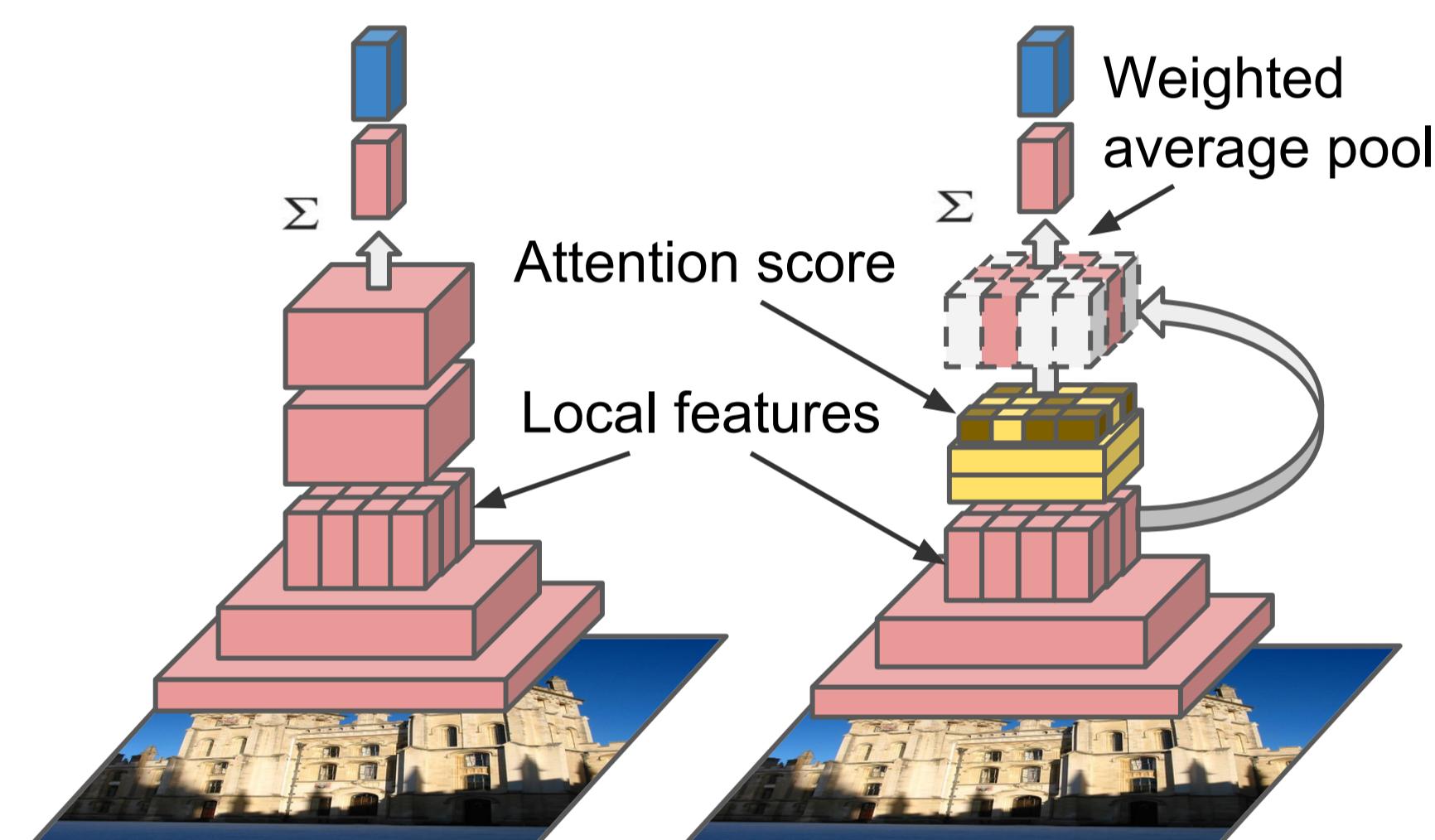
- (1) Descriptor: Fine-tuning CNN based on classification loss
- (2) Attention: Implicit learning by attention-weighted average pooling

## • Weak Supervision

- Image-level annotation from landmarks dataset [Babenko et al., ECCV'14]
- Attention-weighted average pooling

$$y = \mathbf{W} \left( \sum_n \alpha(\mathbf{f}_n; \theta) \cdot \mathbf{f}_n \right)$$

Descriptor learning      Attention learning



## Google-Landmarks Dataset

## □ Construction

- Mined from GPS-tagged photos from the web [Zheng et al., CVPR'09]

## □ Challenges:

- Query with no correct match
- Large / diverse set of landmarks
- Large variations: clutter, occlusion, partially out-of-view object

## □ Diversity of Landmarks / Images



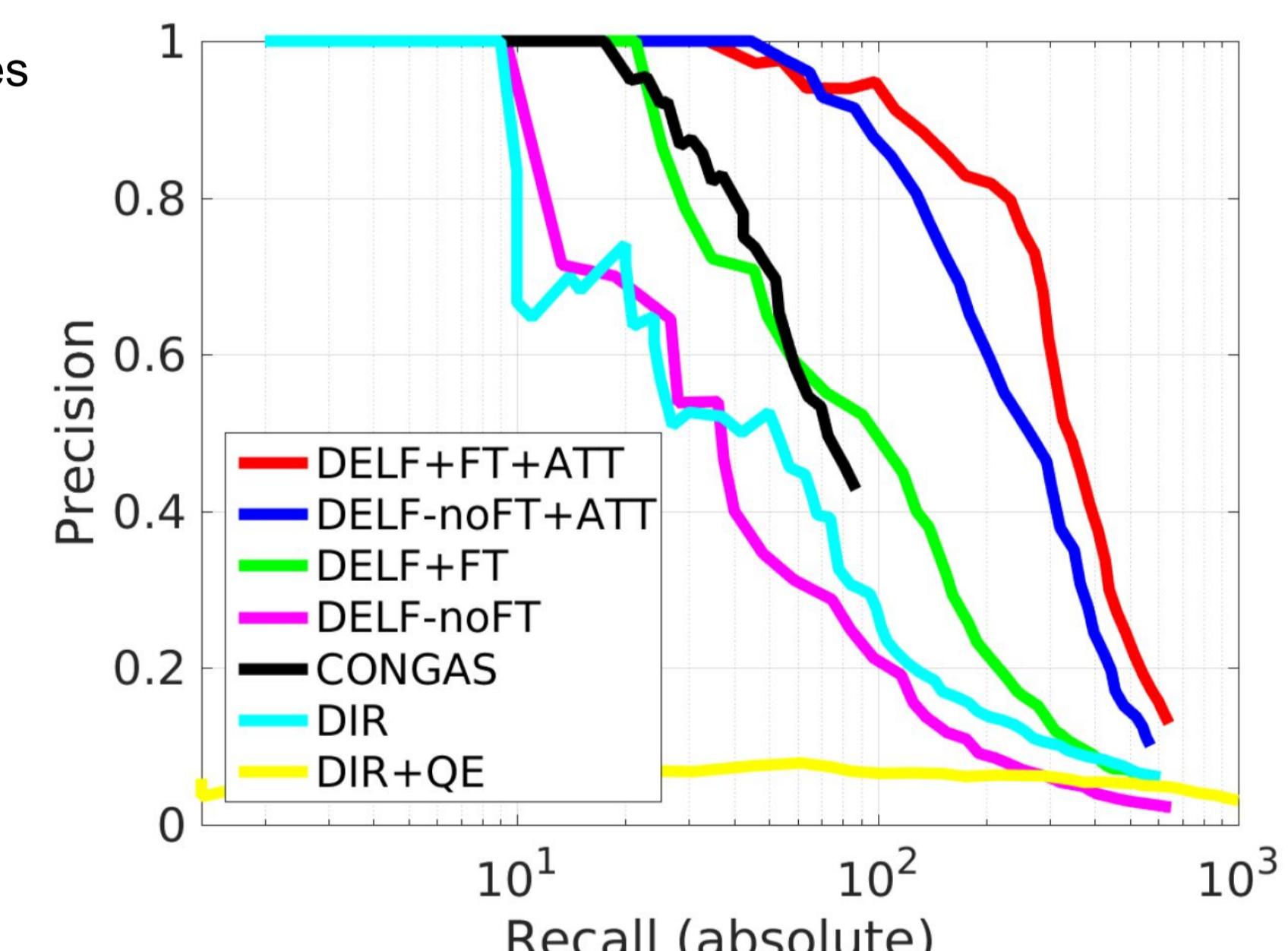
## □ Dataset Benchmark

	# landmarks (scenes)	# query images	# database images	query distractors
Google-Landmarks	12,894	100,000	1,060,709	○
Oxford5k [Philbin et al., CVPR'07]	16	55	5,062	✗
Paris6k [Philbin et al., CVPR'08]	11	55	6,412	✗
Holidays [Jegou et al., ECCV'08]	500	500	1,491	✗

❖ Dataset to be released with Landmark Recognition Challenge

## Experiments

## □ Google-Landmarks



## □ Results on Existing Datasets

Dataset	Oxford5k	Oxford105k	Paris6k	Paris106k
DIR	86.1	82.8	94.5	90.6
DIR+QE	87.1	85.2	95.3	91.8
siaMAC	77.1	69.5	83.9	76.3
siaMAC+QE	81.7	76.6	86.2	79.8
CONGAS	70.8	61.1	67.1	56.8
LIFT	54.0	—	53.6	—
DIR+QE*	89.0	87.8	93.8	90.5
siaMAC+QE*	82.9	77.9	85.6	78.3
DELF+FT+ATT (ours)	83.8	82.6	85.0	81.7
DELF+FT+ATT+DIR+QE (ours)	<b>90.0</b>	<b>88.5</b>	<b>95.7</b>	<b>92.8</b>

- DIR + DELF improves performance significantly
- Complementary information from DELF / DIR

## □ Ablation Study

- Attention helps more than fine-tuning

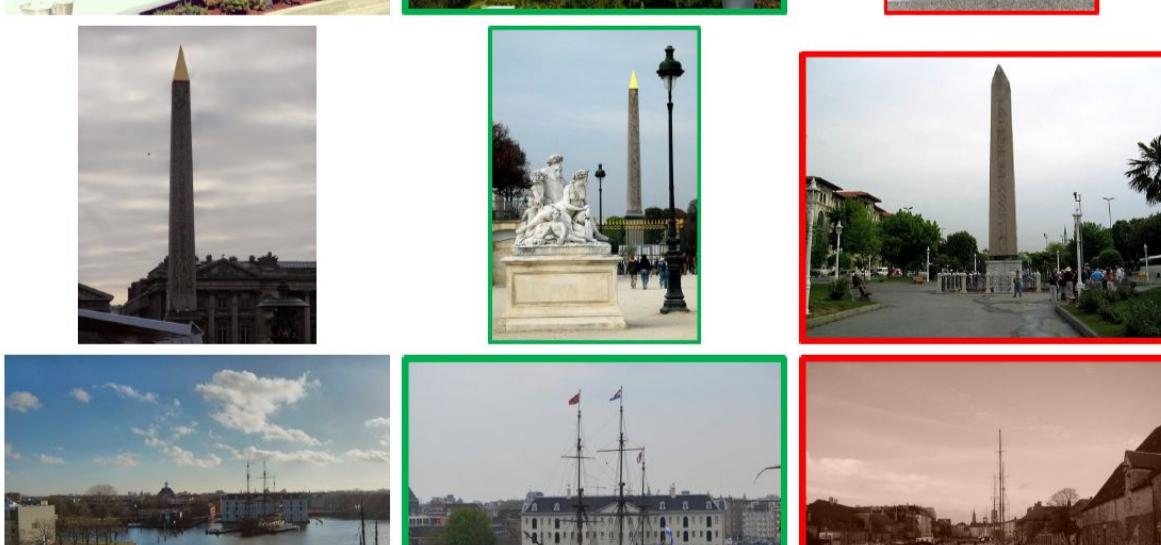
## □ DELF vs DIR

- **DELF**: handles variations better (e.g. scale variations)
- **DIR**: confuses objects with semantic similarity



## □ DELF vs CONGAS

- **DELF** has higher recall



## □ Feature Correspondence

- **CONGAS** fails on following examples: ↓

