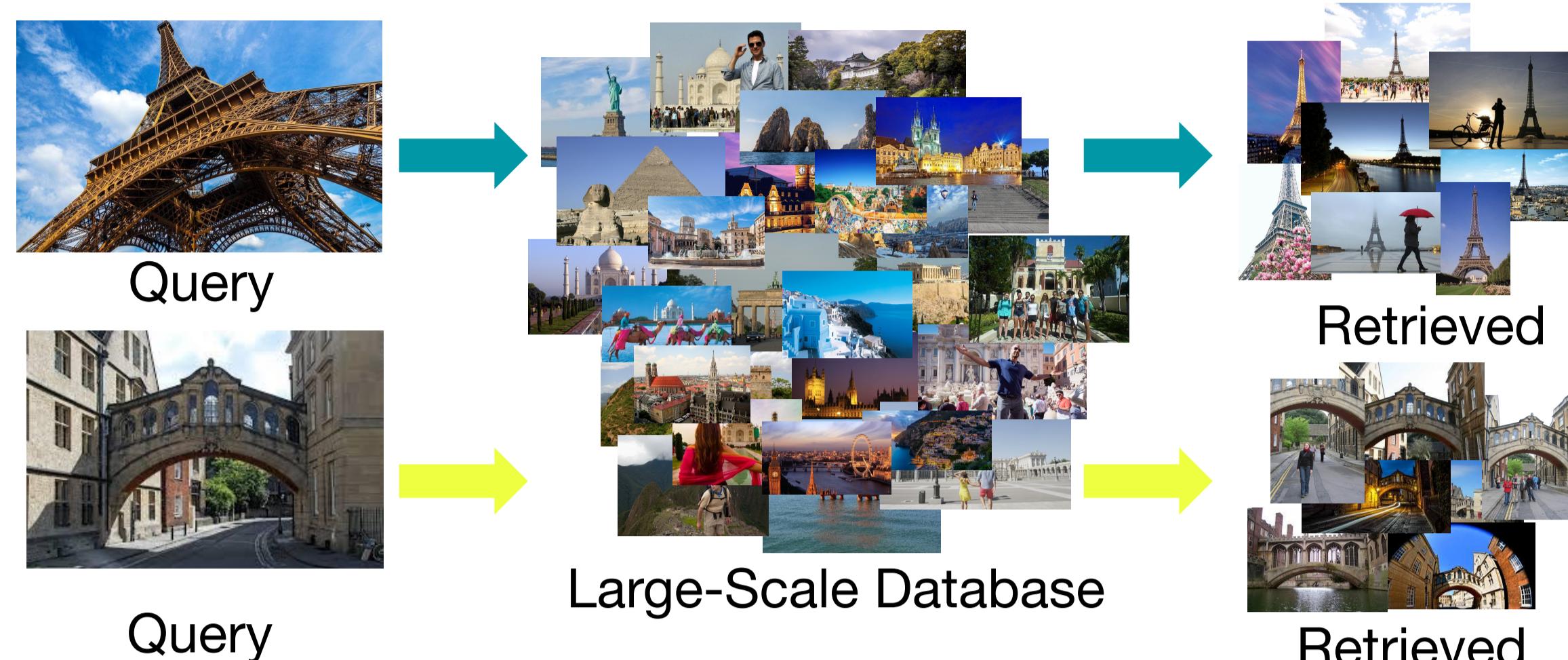


Large-Scale Image Retrieval



Challenges

- Low performance on small objects

- Regional representation using object detector

- Indexing multiple regions is inefficient

- Regional Aggregated Matching Kernels

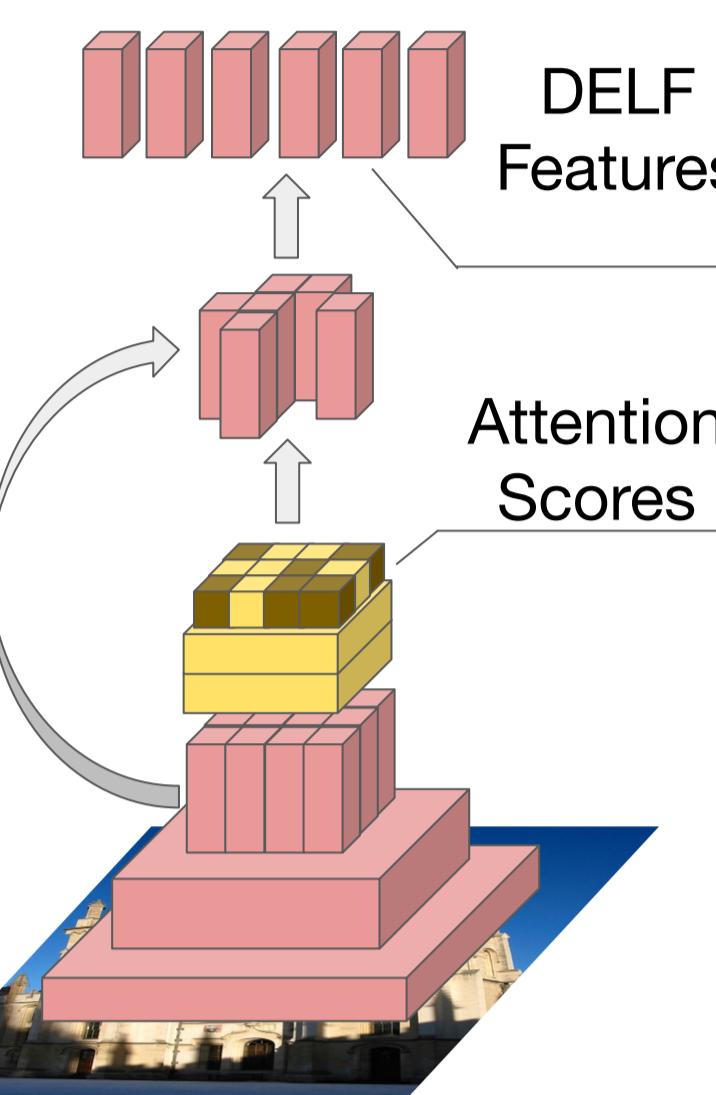
- No manually labeled landmark box datasets

- Google Landmark Boxes dataset

State-of-the-art results on Revisited Oxford/Paris datasets

Contributions

Background: DELF



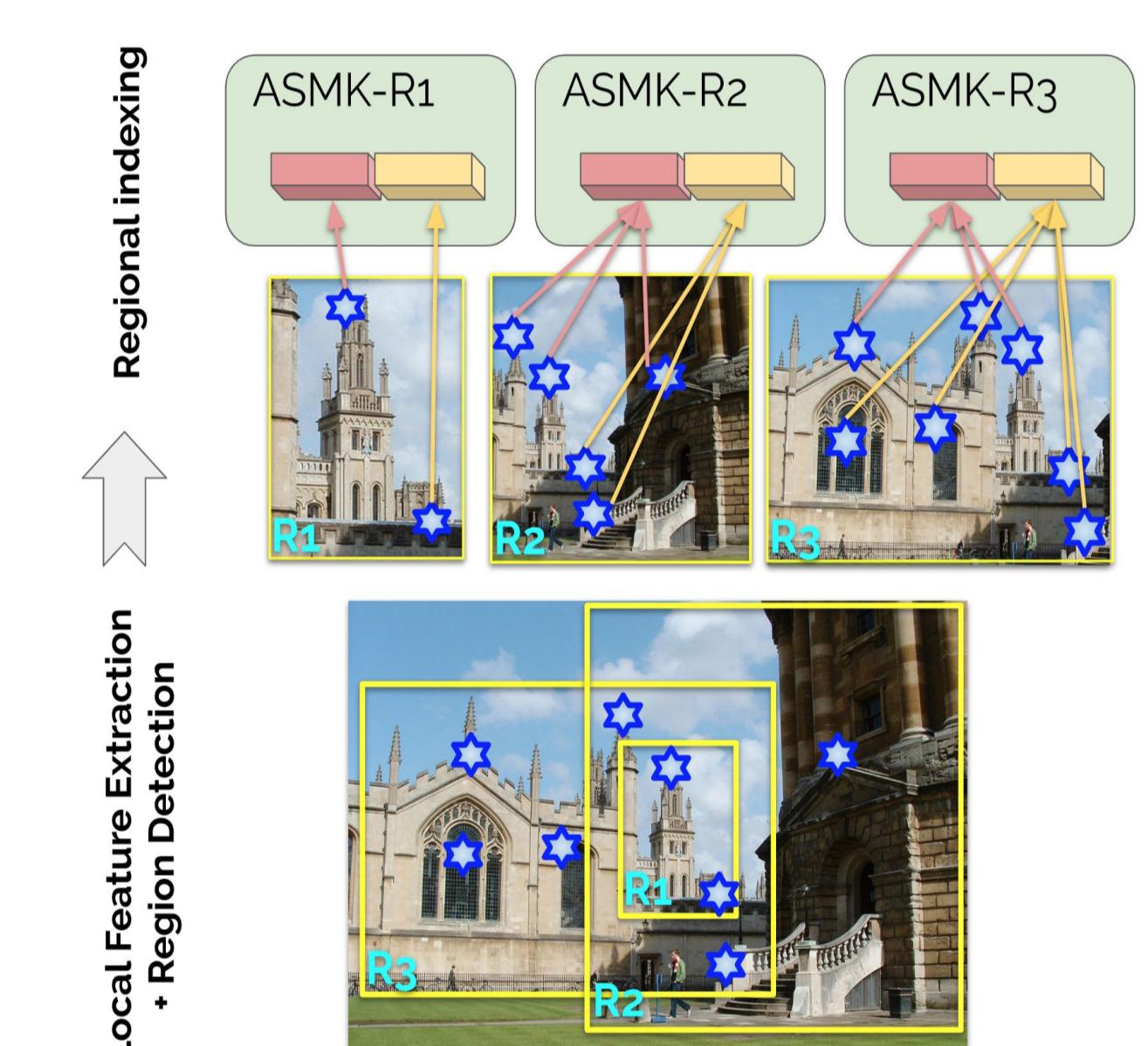
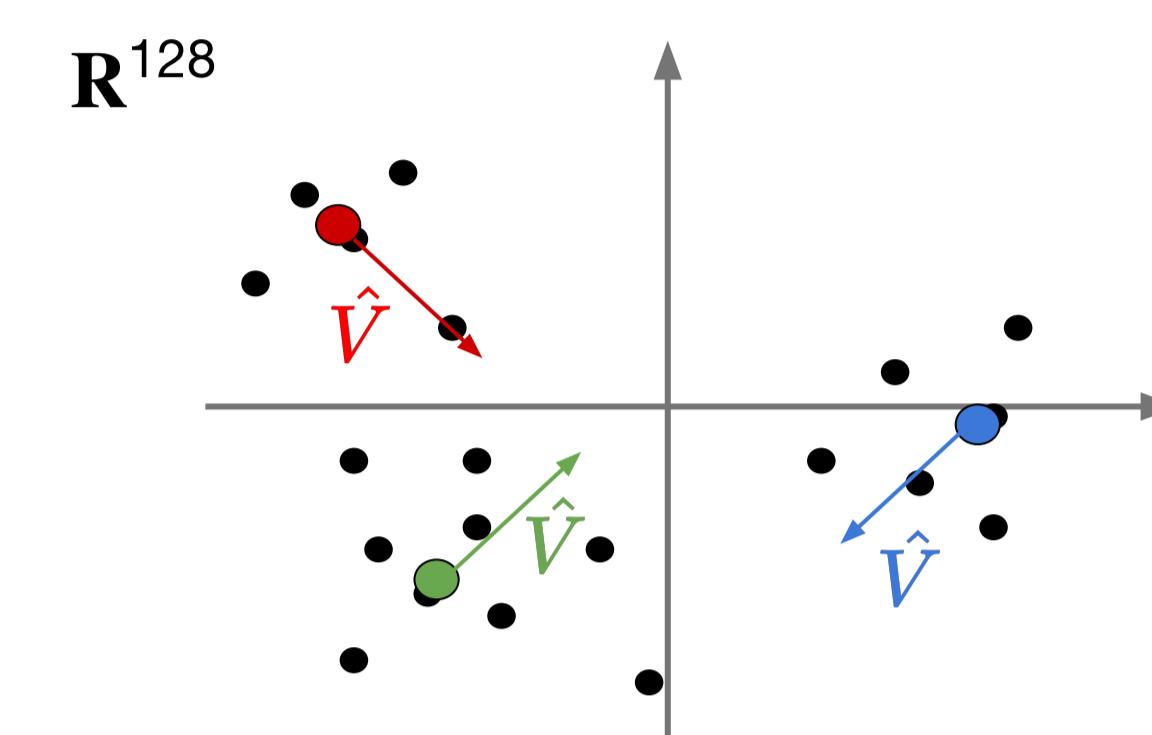
Regional Search

- ✓ DELF-ASMK extracted for each box
- ✓ Each box indexed independently
- ✓ Index size = O(#boxes in all images)
- ✓ Retrieval system:
 - o score for an image = its maximum box score

Background: ASMK

$$\text{sim}^{\text{ASMK}} = \sum_c \text{ASMK}(\mathcal{X}_c, \mathcal{Y}_c)$$

$$\text{ASMK}(\mathcal{X}_c, \mathcal{Y}_c) = \sigma(\hat{V}(\mathcal{X}_c)^\top \hat{V}(\mathcal{Y}_c))$$



Detect-to-Retrieve

Regional Aggregation: R-ASMK

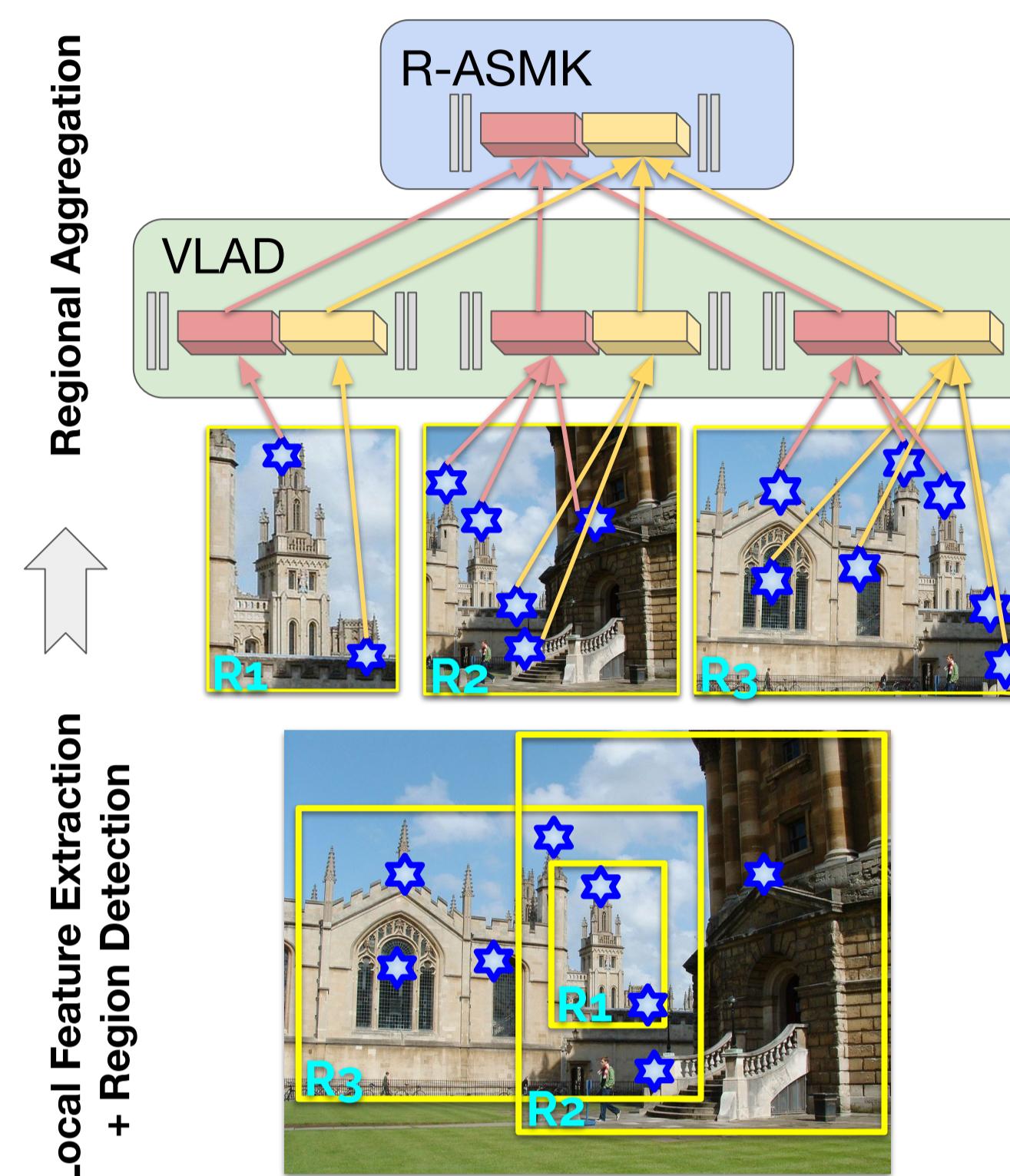
- ✓ Extension of aggregated matching kernels
 - o Pooling over boxes, then images
 - o Boxes drive improved representation
- ✓ Motivation:
 - o average pooling of VLAD similarities
 - o + renormalization
 - o + selectivity function
- ✓ Index size = O(#images)
 - Memory-efficient, especially important for large-scale setting

$$\text{sim}^{\text{R-ASMK}} = \sum_c \text{R-ASMK}(\{\mathcal{X}_c^r\}_r, \{\mathcal{Y}_c^r\}_r)$$

$$\text{R-ASMK}(\{\mathcal{X}_c^r\}_r, \{\mathcal{Y}_c^r\}_r) = \sigma(\hat{V}_R(\{\mathcal{X}_c^r\}_r)^\top \hat{V}_R(\{\mathcal{Y}_c^r\}_r))$$

$$\hat{V}_R(\{\mathcal{Y}_c^r\}_r) = \frac{V_R(\{\mathcal{Y}_c^r\}_r)}{\|V_R(\{\mathcal{Y}_c^r\}_r)\|}$$

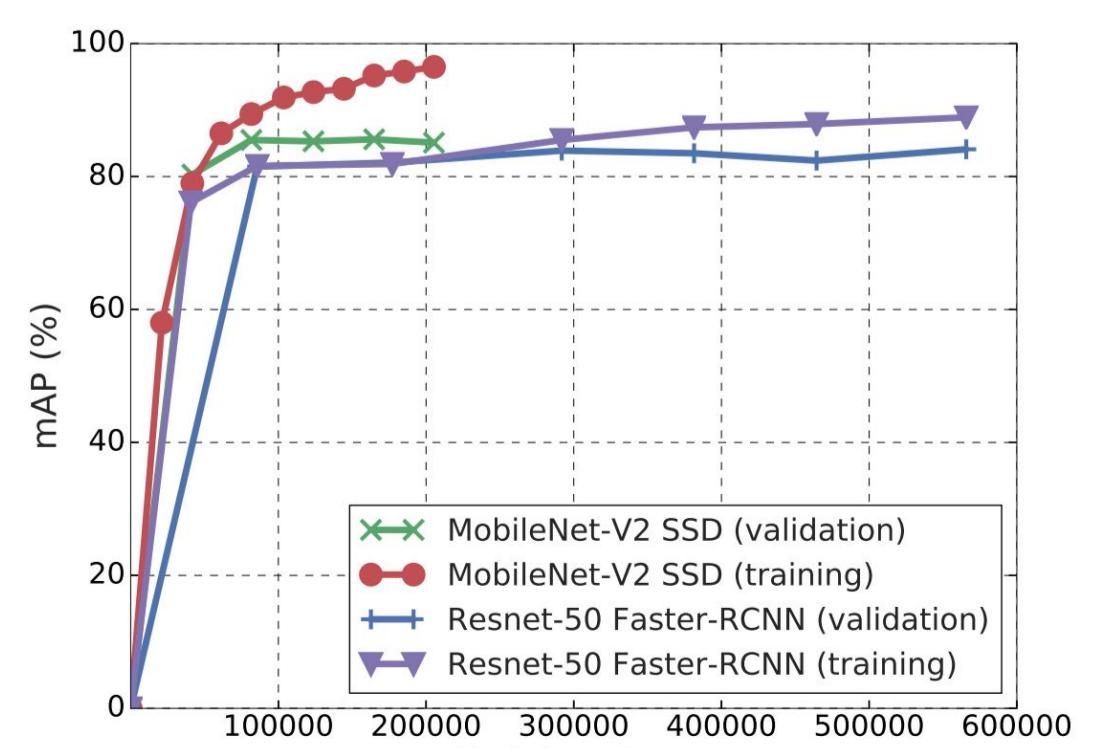
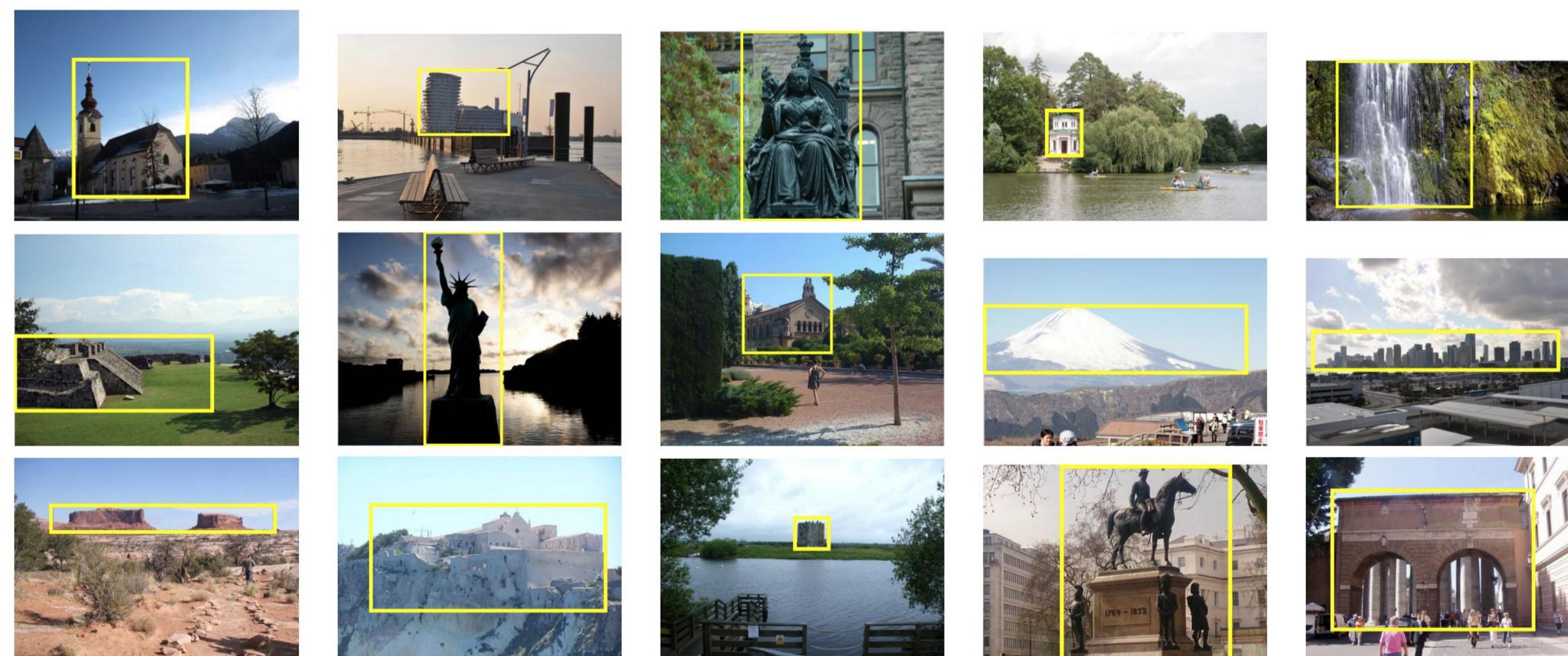
$$V_R(\{\mathcal{Y}_c^r\}_r) = \frac{1}{R_n} \sum_r \gamma(\mathcal{Y}_c^r) V(\mathcal{Y}_c^r)$$



- ✓ Alternative naive regional kernels
 - o E.g., R-VLAD: average pooling of VLAD's
 - o Normalization issues: per-VW information is "watered down"

Google Landmark Boxes

- ✓ 86k annotated boxes, from 15k landmarks
- ✓ One box per image capturing most prominent landmark
- ✓ Accurate detection with off-the-shelf architectures

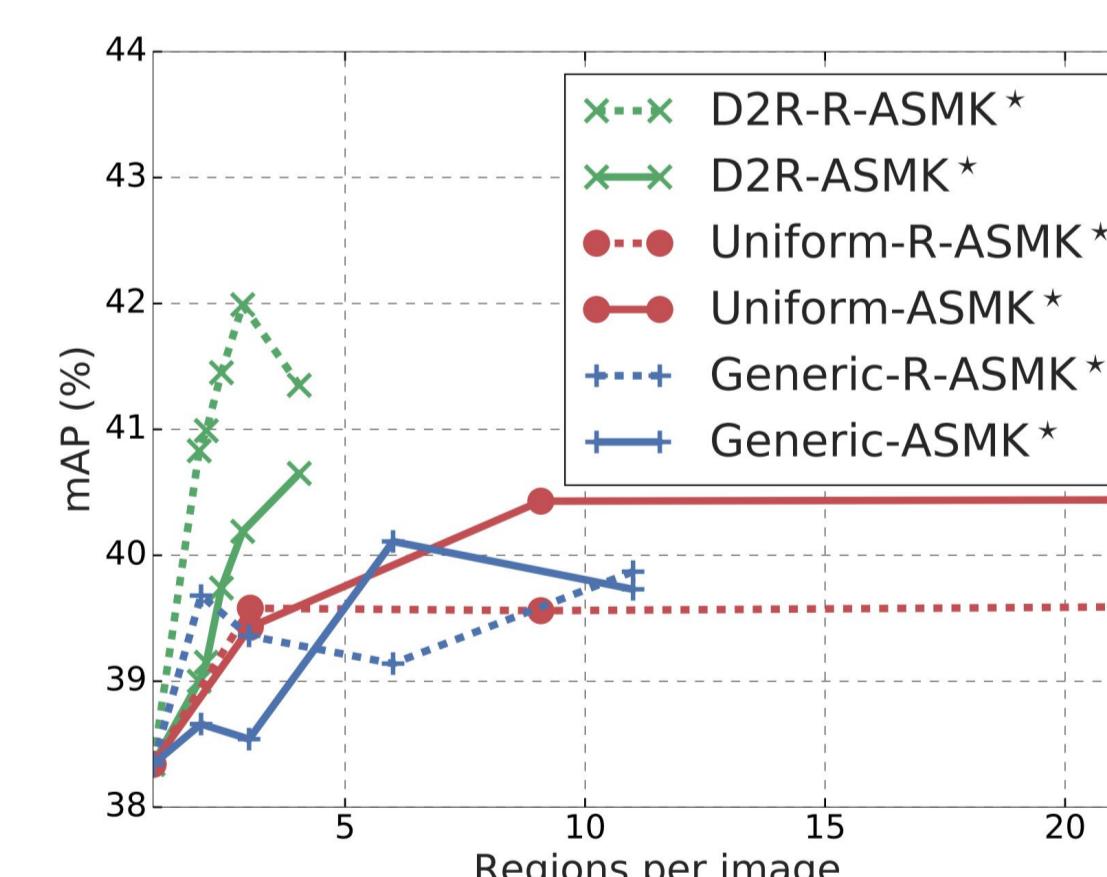


<https://www.kaggle.com/google/google-landmarks-dataset>

Experimental Results

Ablation Study

- ✓ D2R improves upon no-detection baseline
 - o 2.31% for regional search
 - o 3.65% for regional aggregation
- ✓ Regional aggregation > Regional Search
 - o Higher mAP, smaller index
- ✓ D2R > uniform, generic detectors



State-of-the-Art Image Retrieval Results

- ✓ Best results in Revisited datasets for all protocols and metrics
- ✓ 9.3% mAP improvement on Revisited Oxford (Hard)
- ✓ 1.9% mAP improvement on Revisited Paris (Hard)

Method	Medium				Hard			
	\mathcal{R}_{Oxf} mAP	$\mathcal{R}_{Oxf} + \mathcal{R}_{IM}$ mAP	\mathcal{R}_{Par} mAP	$\mathcal{R}_{Par} + \mathcal{R}_{IM}$ mAP	\mathcal{R}_{Oxf} mAP	$\mathcal{R}_{Oxf} + \mathcal{R}_{IM}$ mAP	\mathcal{R}_{Par} mAP	$\mathcal{R}_{Par} + \mathcal{R}_{IM}$ mAP
AlexNet-GeM [30]	43.3	62.1	24.2	42.8	58.0	91.6	29.9	84.6
VGG16-GeM [30]	61.9	82.7	42.6	68.1	69.3	97.9	45.4	94.1
ResNet101-R-MAC [9]	60.9	78.1	39.3	62.1	78.9	96.9	54.8	93.9
ResNet101-GeM [30]	64.7	84.7	45.2	71.7	77.2	98.1	52.3	95.3
ResNet101-GeM+DSM [34]	65.3	87.4	47.6	76.4	77.4	99.1	52.8	96.7
HeSift-rSIFT-ASMK* [38]	60.4	85.6	45.0	76.0	61.2	97.9	42.0	95.3
HeSift-rSIFT-ASMK*+SP [38]	60.6	86.1	46.8	79.6	61.4	97.9	42.3	95.3
HeSift-HardNet-ASMK*+SP [24]	65.6	90.2	49.2	82.6	65.2	98.9	41.1	97.9
DELF-ASMK*+SP [25, 28]	67.8	87.9	53.8	81.1	76.9	99.3	57.3	98.3
DELF-ASMK* (reimpl.)	65.7	87.9	—	—	77.1	98.7	—	—
DELF-D2R-R-ASMK* (ours)	69.9	89.0	—	—	78.7	99.0	—	—
— DELF-GLD (ours)	73.3	90.0	61.0	84.6	80.7	99.1	60.2	97.9
DELF-ASMK*+SP (reimpl.)	68.9	90.9	—	—	76.6	98.7	—	—
DELF-D2R-R-ASMK*+SP (ours)	71.9	91.3	—	—	78.0	99.4	—	—
— DELF-GLD (ours)	76.0	93.4	64.0	87.7	80.2	99.1	59.7	99.0
DELF-ASMK*+SP (reimpl.)	68.9	90.9	—	—	76.6	98.7	—	—
DELF-D2R-R-ASMK*+SP (ours)	71.9	91.3	—	—	78.0	99.4	—	—
— DELF-GLD (ours)	76.0	93.4	64.0	87.7	80.2	99.1	59.7	99.0

