

Class-Weighted Convolutional Features for Image Retrieval



[Albert Jiménez](#)



[Xavier Giró-i-Nieto](#)



[Jose Alvarez](#)



[\[arXiv\]](#)



[\[Code\]](#)



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Image Retrieval

Google Search By Image

The screenshot shows a Google search results page for an image of a modern building. At the top, there is a large black rectangular placeholder for the image. Below it, the search query is "Best guess for this image: universitat politècnica de catalunya barcelona".

The first result is a link to "UPC. Universitat Politècnica de Catalunya, BarcelonaTech" (www.upc.edu/estat_lenguajes). The snippet describes UPC as a public university dedicated to engineering and architecture.

The second result is a link to "Polytechnic University of Catalonia - Wikipedia" (https://en.wikipedia.org/w/index.php?title=Polytechnic_University_of_Catalonia&oldid=8001000). The snippet provides a brief history of the university's formation in 1971 through the merger of engineering and architecture schools.

Below the search results, there is a section titled "Visually similar images" which displays a grid of nine smaller images showing various modern buildings.

On the right side of the page, there is a sidebar with information about the Polytechnic University of Catalonia. It includes:

- Basic info: Province of Barcelona, Total enrollment: 36,000 (2010), Founded: March 1971, Endowment: 263.5 million EUR.
- Social media links for Twitter, YouTube, and Facebook.
- A "Notable alumni" section featuring four profile pictures: Jordi Guitart, Eduard Vayreda, Carme Torras, and Antoni Lluis-Masó.
- A "People also search for" section with links to UOC, Universitat Pompeu Fabra, and UPC.

Image Retrieval

Task Definition

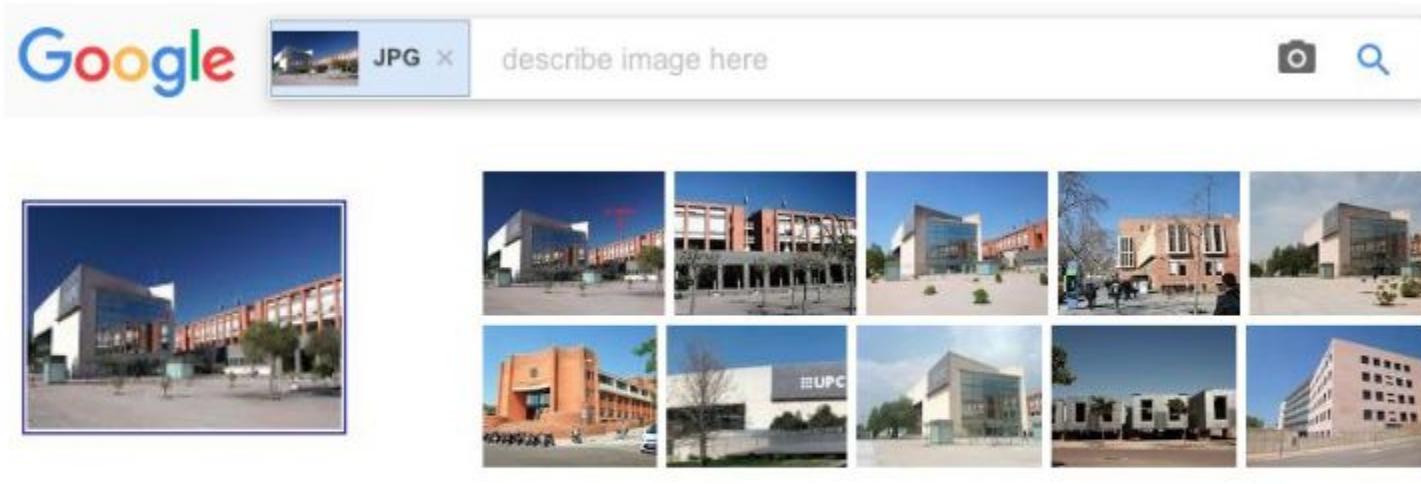


Image Query

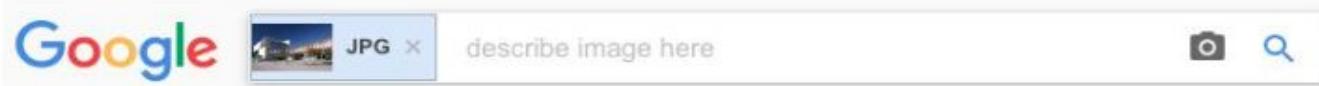
Visual Similar Images

Best guess for this image: *universitat politècnica de catalunya*

Given an image query, generate a ranked list of similar images

Image Retrieval

Task Definition



The image shows a screenshot of a Google Image search results page. At the top left is the Google logo. Next to it is a small thumbnail image labeled "JPG" with a close button "X". To the right of the thumbnail is a text input field containing "describe image here". Further to the right are a camera icon and a magnifying glass search icon. Below the search bar is a grid of images. On the left, there is a large thumbnail labeled "Image Query" below it. To the right of the query is a grid of smaller images labeled "Visual Similar Images" below it.

Image Query

Visual Similar Images

Best guess for this image: *universitat politècnica de catalunya*

Given an image query, generate a ranked list of similar images

What are the properties that we want for these systems?

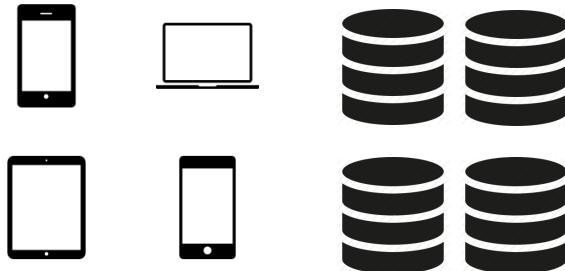
Image Retrieval

Desirable System Properties

Invariant to scale, illumination, translation



Memory Efficient



Be able to provide a fast search

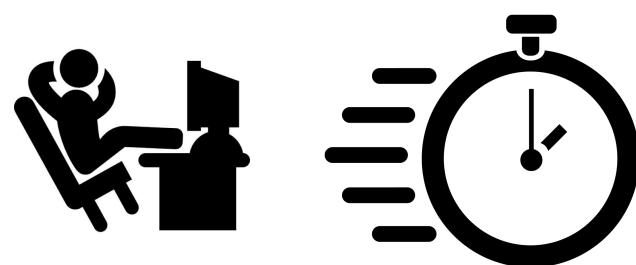


Image Retrieval

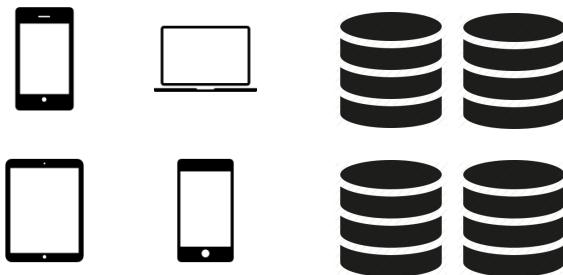
Desirable System Properties

Invariant to scale, illumination, translation

How to achieve them?



Memory Efficient



Be able to provide a fast search

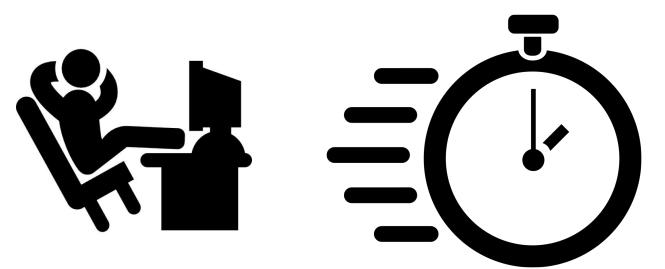
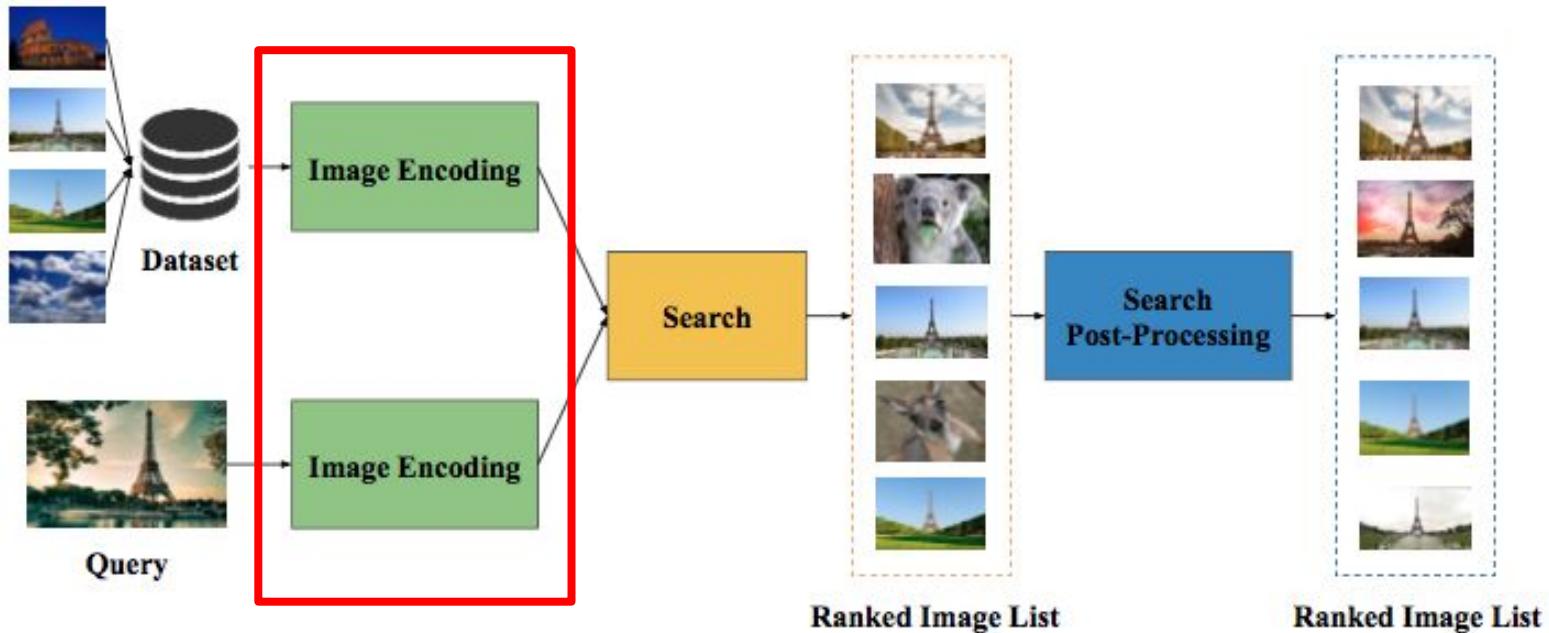


Image Retrieval

General Pipeline

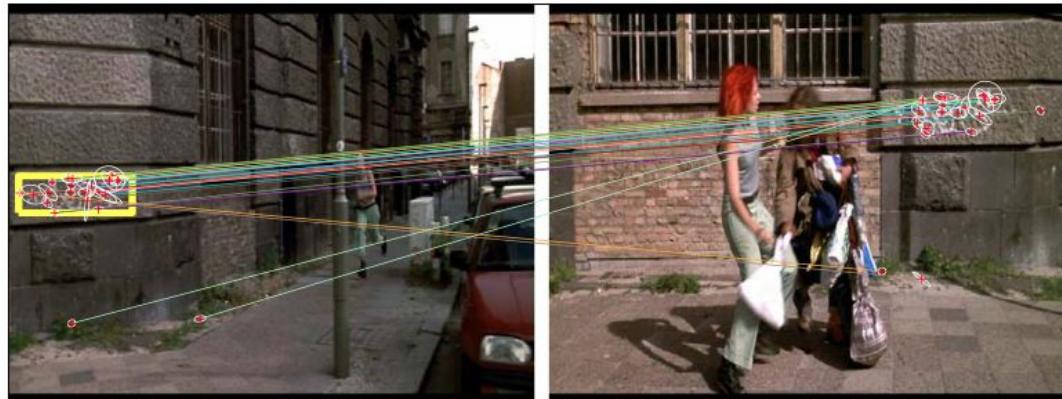


- Encode images into representations that are
 - Invariant
 - Compact

Image Retrieval

Image Representations

- Hand-Crafted Features (Before 2012)
 - SIFT
 - SIFT + BoW
 - VLAD
 - Fisher Vectors

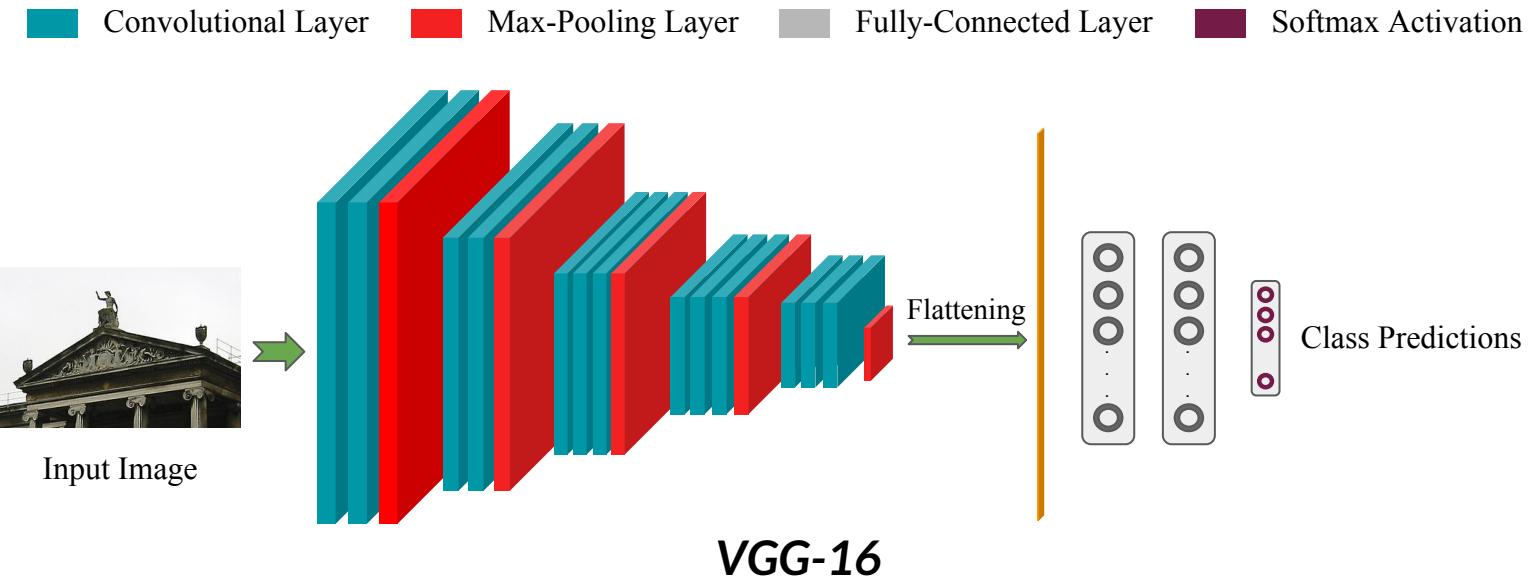


Josef Sivic, Andrew Zisserman, et al. [Video google: A text retrieval approach to object matching in videos](#). In ICCV, 2003.

Image Retrieval

Image Representations

- Learned Features (After 2012)
 - Convolutional Neural Networks (CNNs)



Karen Simonyan and Andrew Zisserman. **Very deep convolutional networks for large-scale image recognition.** *arXiv preprint arXiv:1409.1556.* (2014)

Image Retrieval

Image Representations

- Learned Features (After 2012)
 - Convolutional Neural Networks (CNNs)

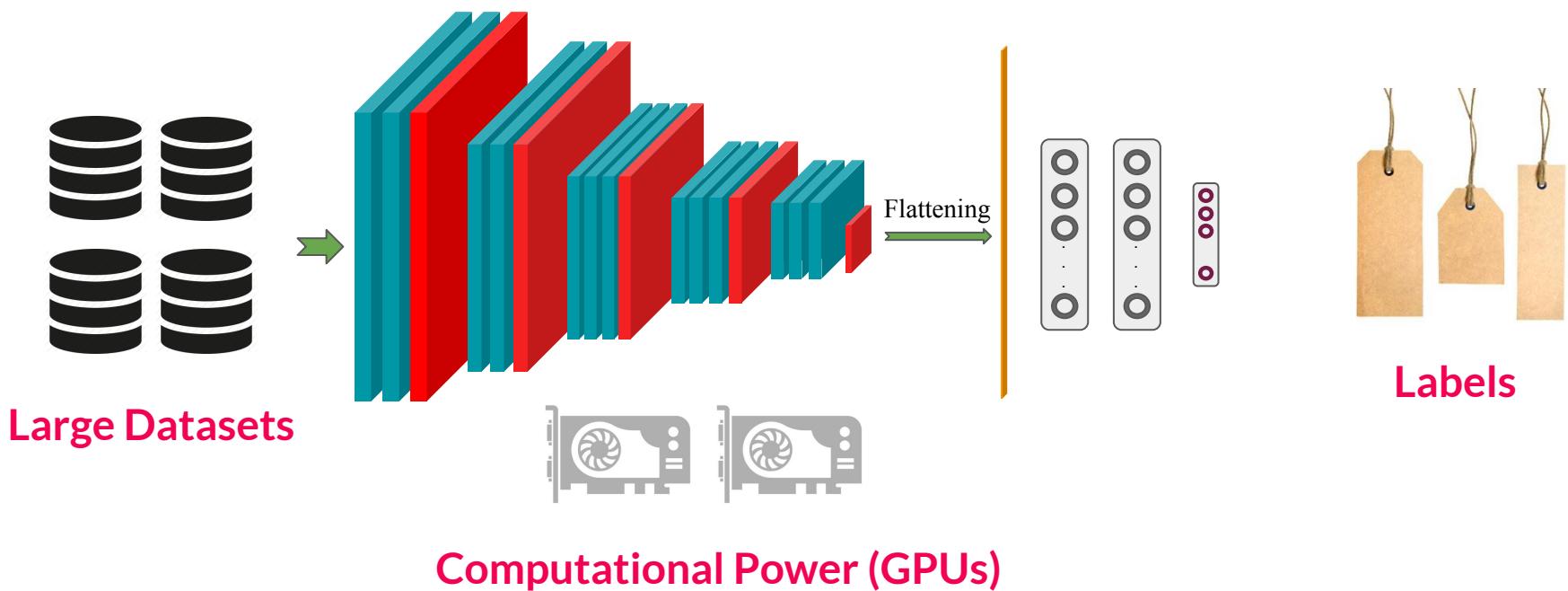
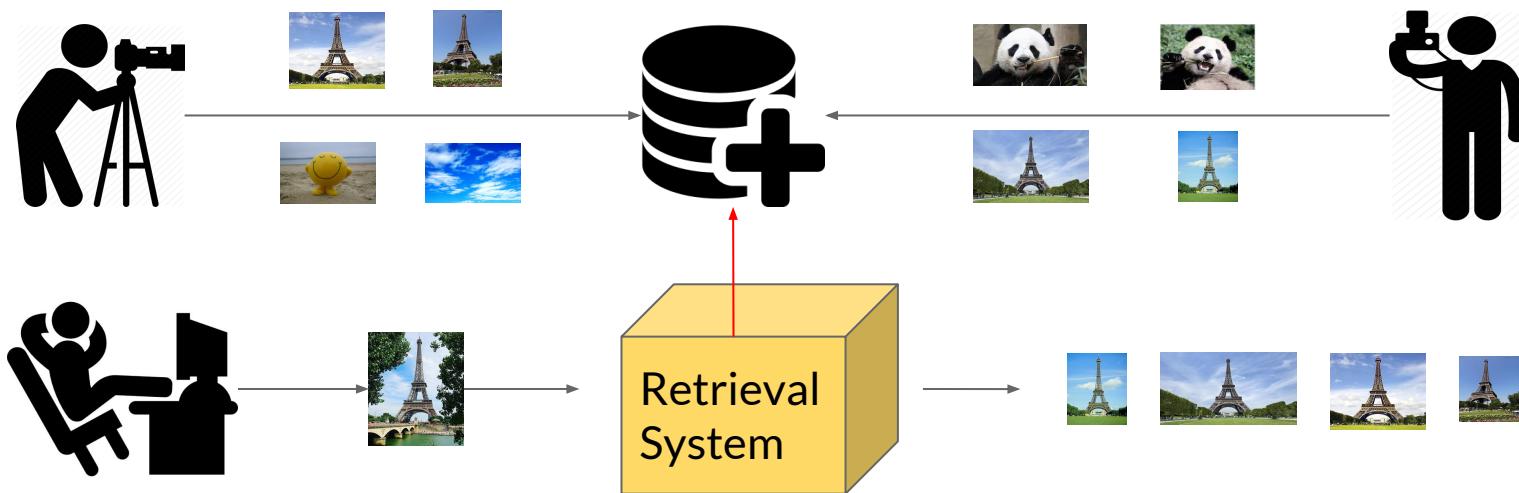


Image Representations

Challenges for Learning Features in Retrieval

Dynamic datasets of unlabeled images



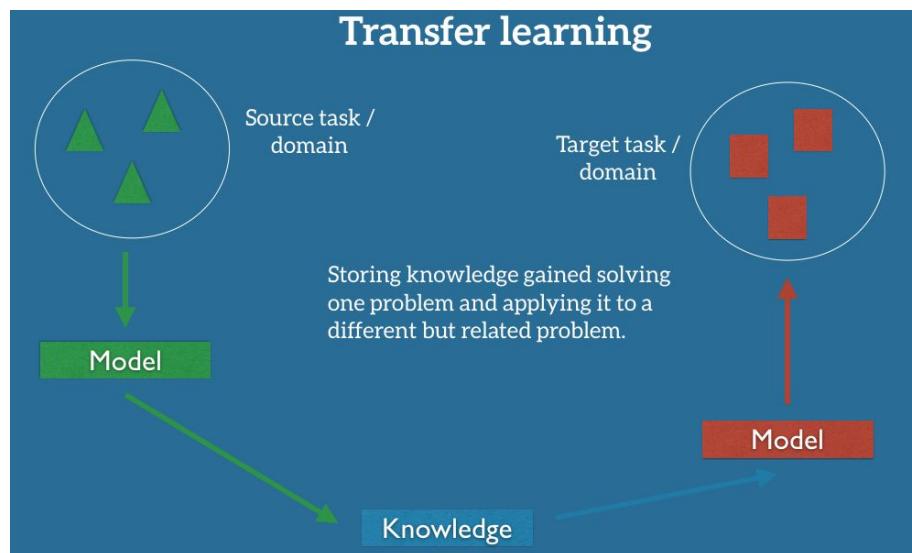
- Collecting, annotating and cleaning a large dataset to train models requires great effort.
- Training or fine-tuning a model every time new data is added is neither efficient nor scalable.

Image Representations

Challenges for Learning Features

One solution that does not involve training is **Transfer Learning**

- Transferring the knowledge of a model trained in a large-scale dataset.
(e.g. ImageNet)



Classification → Retrieval

How?

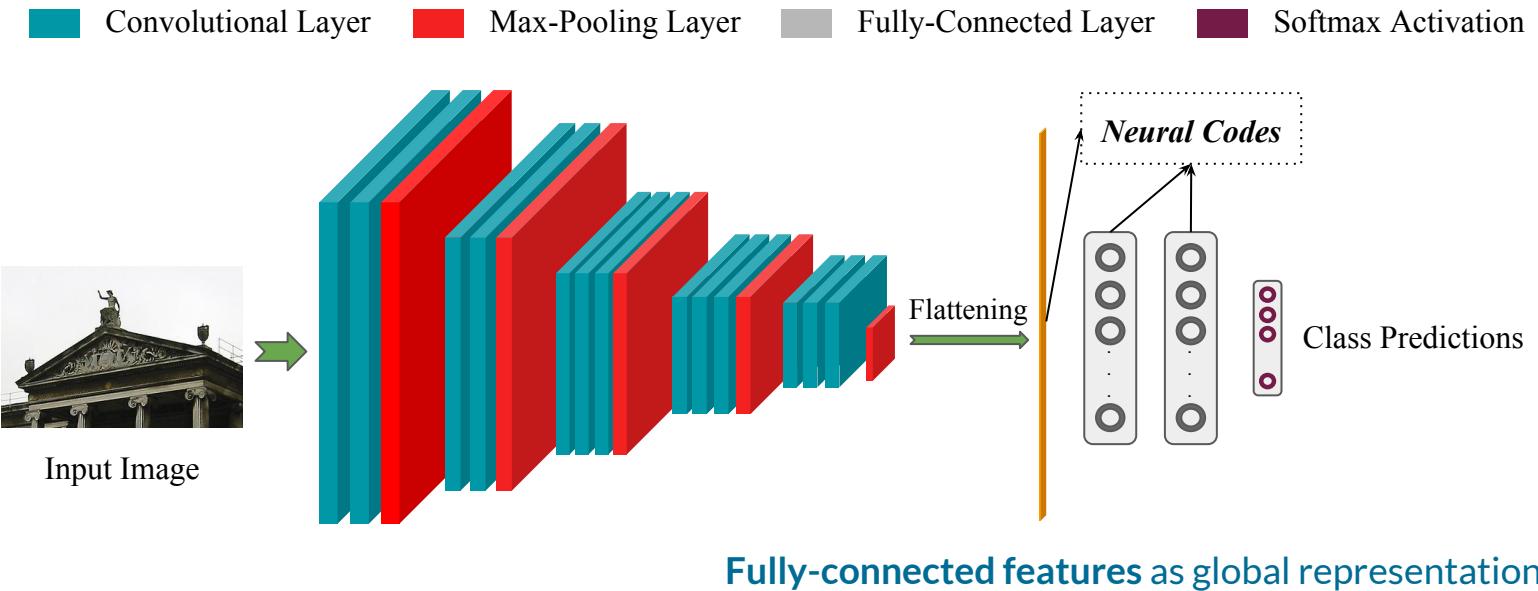
Outline

- ▷ Related Work
- ▷ Proposal
- ▷ Experiments
- ▷ Conclusions

Related Work

Image Representations

Fully-Connected Features

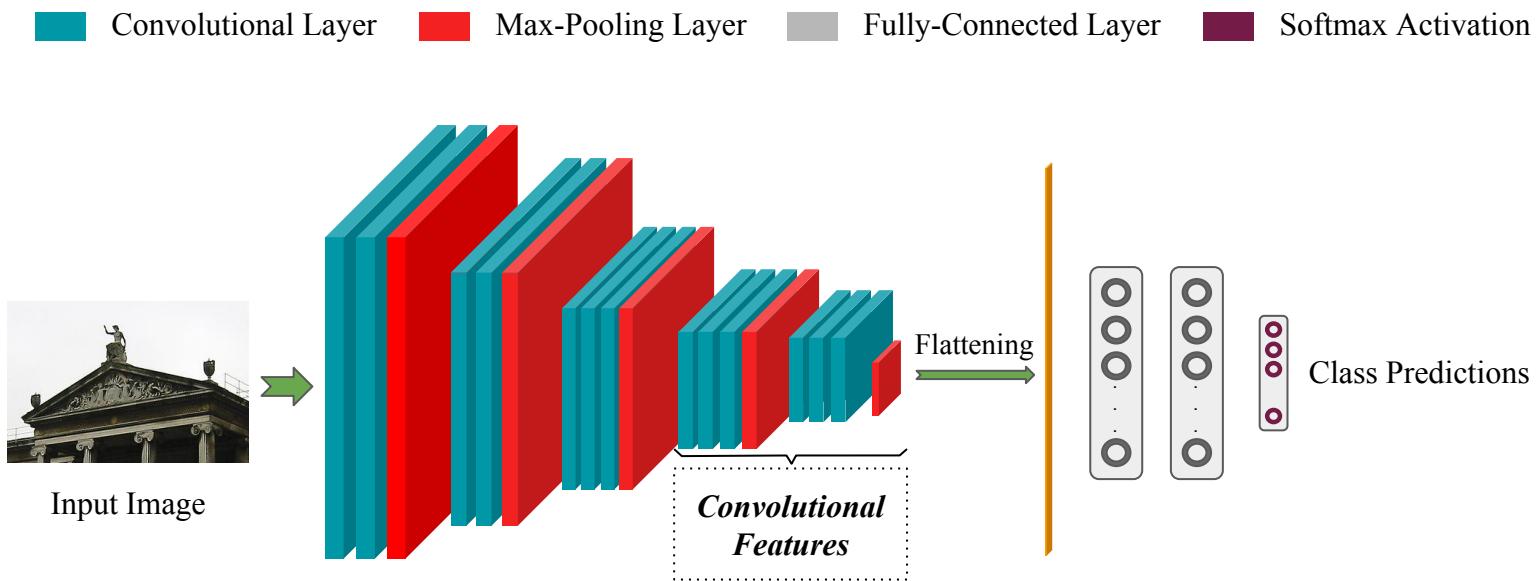


Babenko, A., Slesarev, A., Chigorin, A., & Lempitsky, V. (2014). **Neural codes for image retrieval**. In ECCV 2014

Razavian, A., Azizpour, H., Sullivan, J., & Carlsson, S. (2014). **CNN features off-the-shelf: an astounding baseline for recognition**. In DeepVision CVPRW 2014

Image Representations

Convolutional Features



Convolutional features as global representation
→ They convey the **spatial information** of the image

Babenko, A., & Lempitsky, V. (2015). Aggregating local deep features for image retrieval. *ICCV 2015*

Tolias, G., Sicre, R., & Jégou, H. (2016). Particular object retrieval with integral max-pooling of CNN activations. *ICLR 2016*

Kalantidis, Y., Mellina, C., & Osindero, S. (2015). Cross-dimensional Weighting for Aggregated Deep Convolutional Features. *arXiv preprint arXiv:1512.04065*.

Image Representations

From Feature Maps to Compact Representations

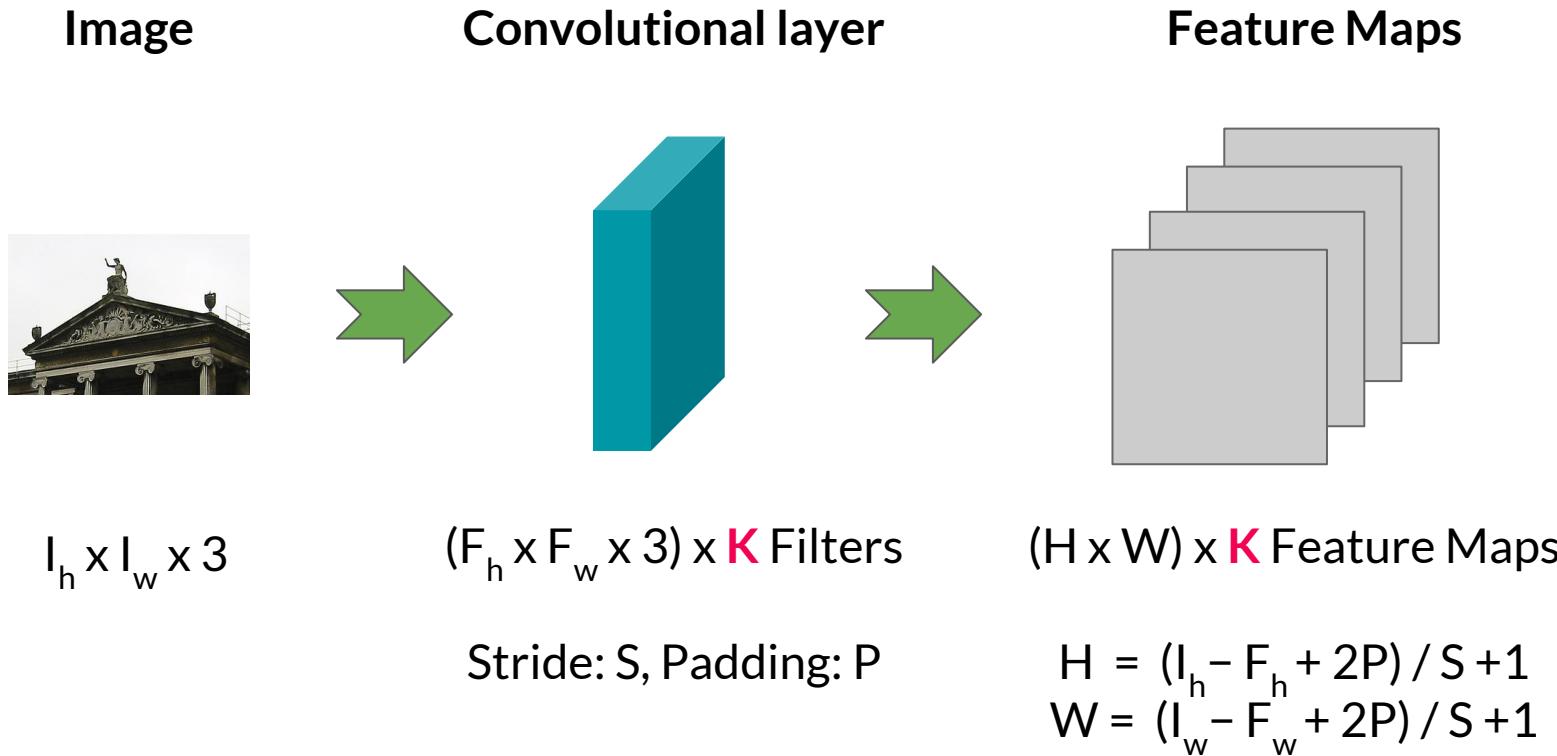


Image Representations

From Feature Maps to Compact Representations

3 Feature Maps (3x3)

1	0	1
0	1	0
9	0	4

Pooling Operation

Max Pooling

9

1	2	1
0	0	0
0	0	0

Max Pooling

2

0	0	1
0	0	0
5	1	1

Max Pooling

5

Image Representations

From Feature Maps to Compact Representations

3 Feature Maps (3x3)

1	0	1
1	1	0
9	0	4

Pooling Operation

Sum Pooling

17

1	2	1
0	0	0
0	0	0

Sum Pooling

4

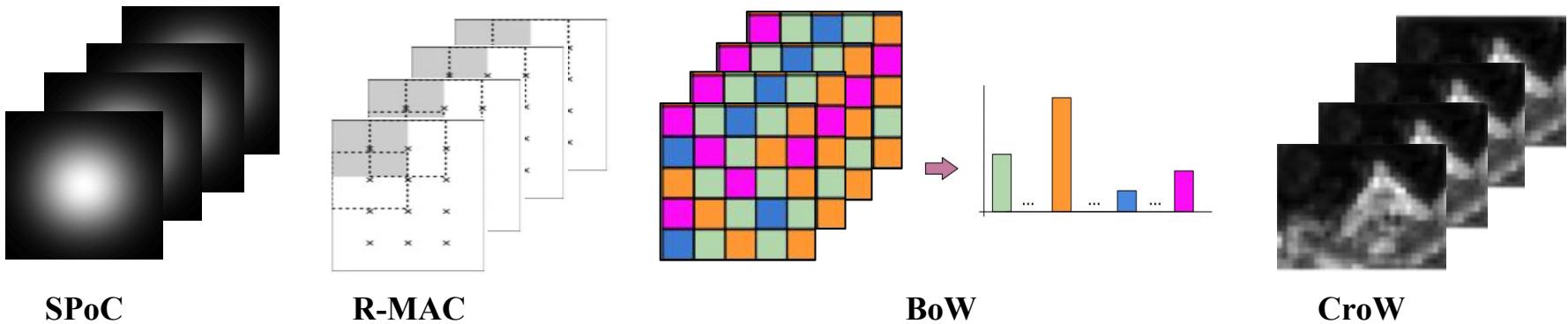
0	0	1
0	0	0
5	1	1

Sum Pooling

8

Image Representations

Related Works

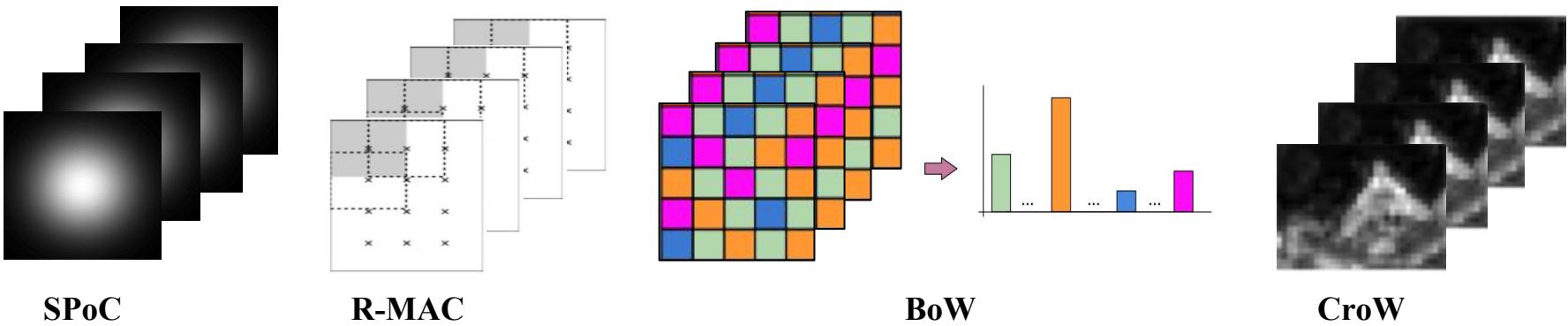


Combine and aggregate convolutional features:

- **SPoC** → Gaussian prior + Sum Pooling
- **R-MAC** → Rigid grid of regions + Max Pooling
- **BoW** → Bag of Visual Words
- **CroW** → Spatial weighting based on normalized sum of feature maps (boost salient content) + channel weighting based on sparsity + Sum Pooling

Image Representations

Related Works



All except CroW:

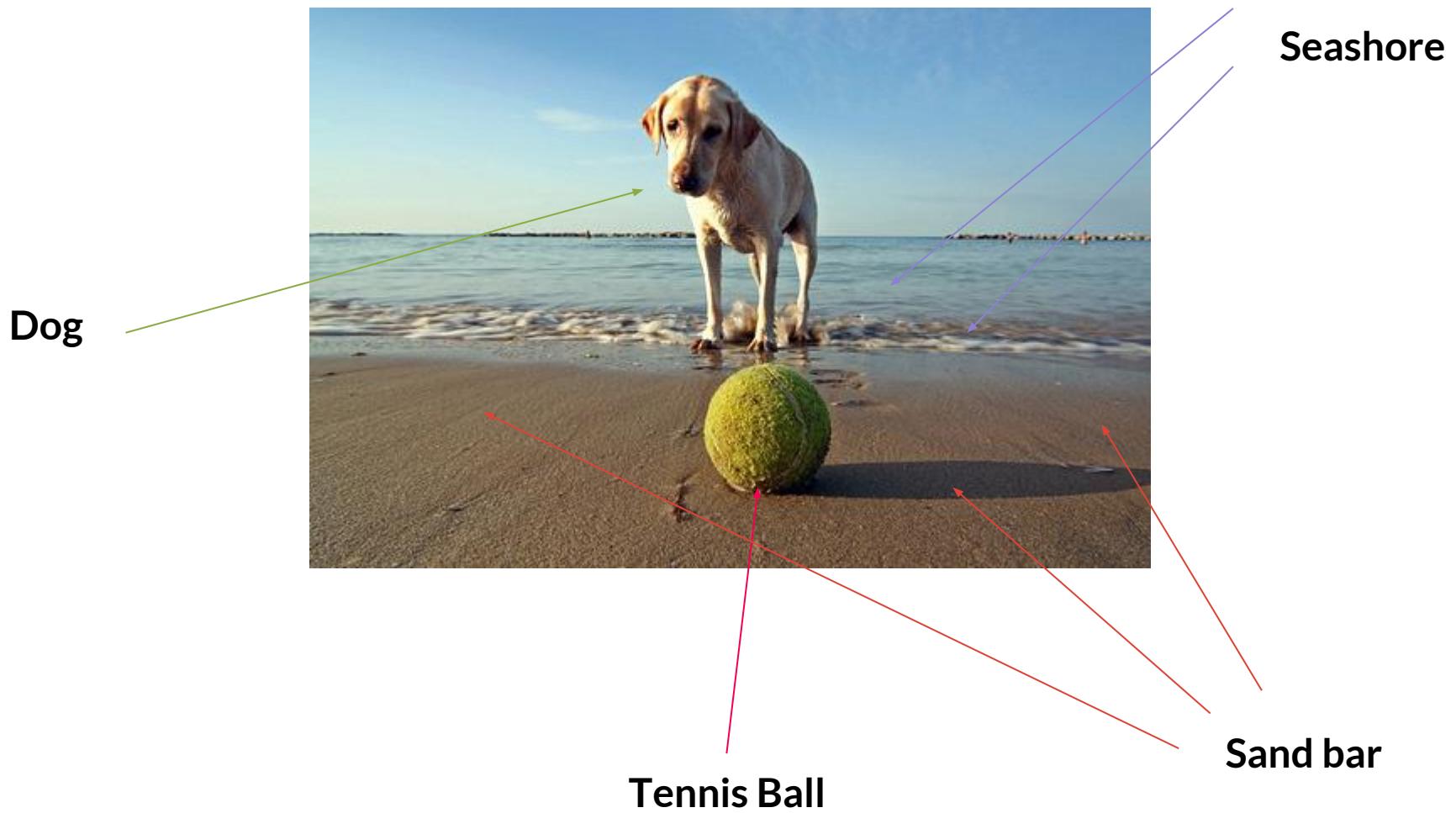
- Apply fixed weights
- Based on fixed regions

Not based on image contents

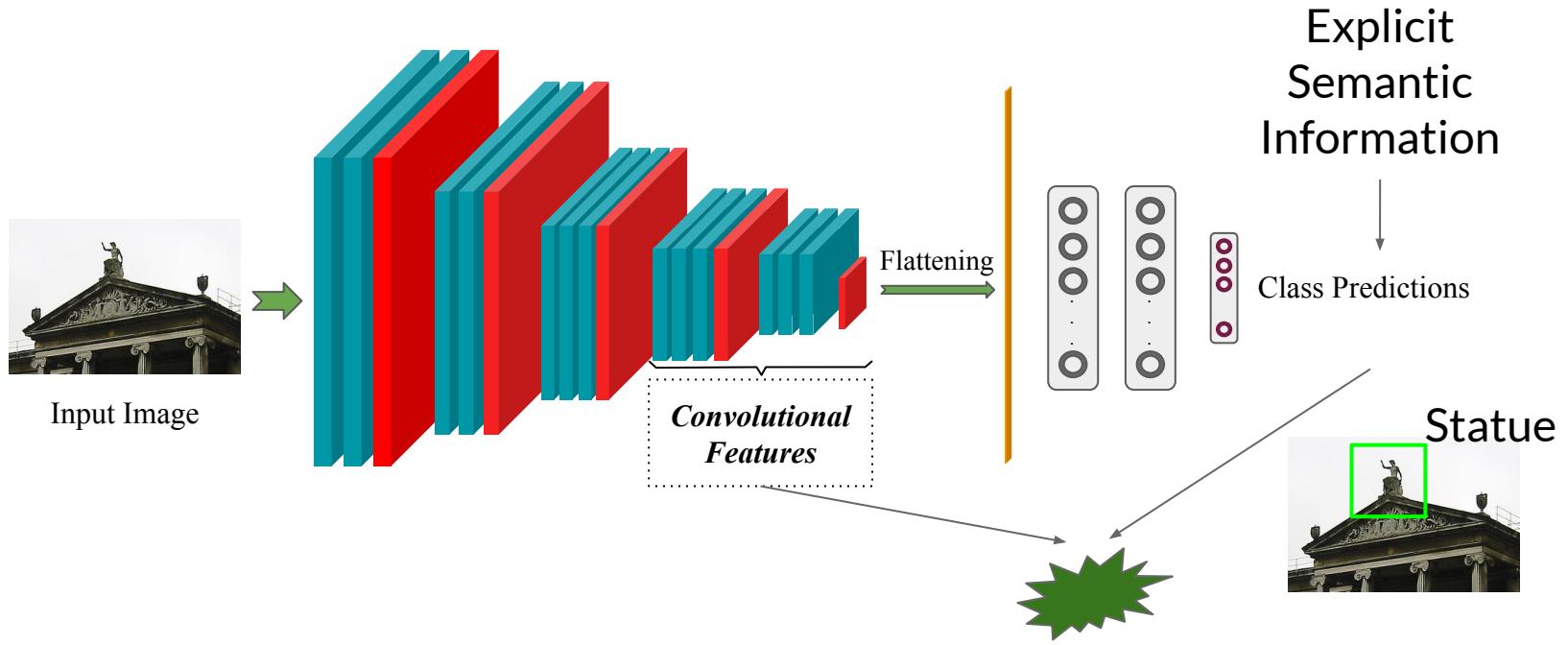
Proposal

Motivation

Image Semantics



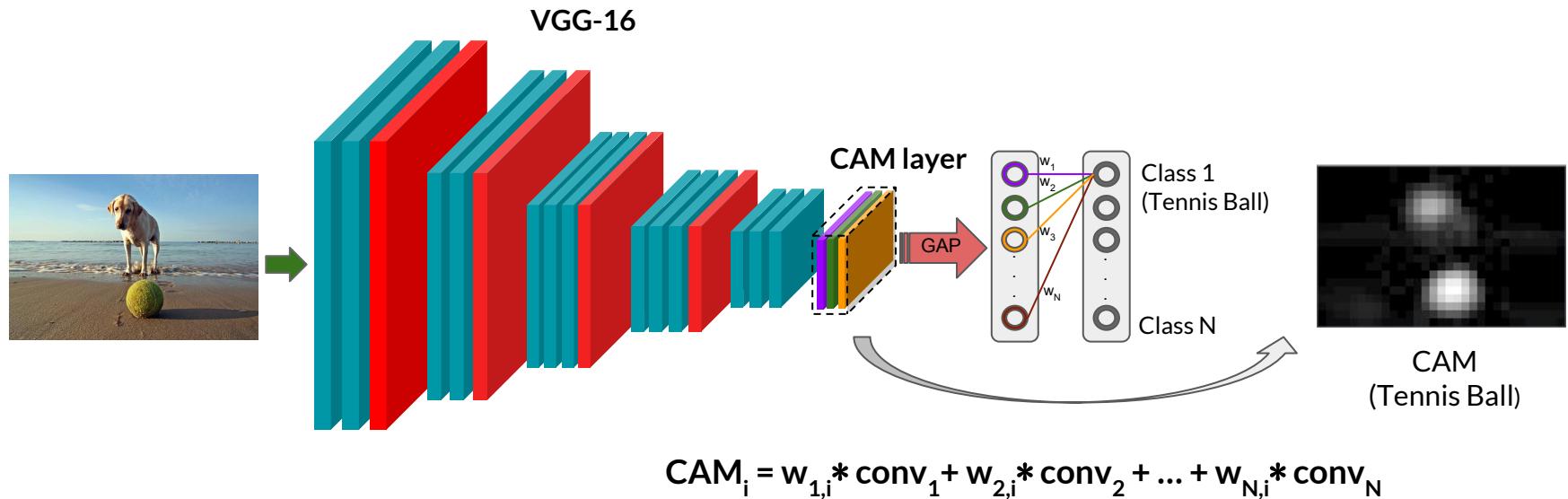
Objective



Combine semantic content of the images with convolutional features

Class Activation Maps

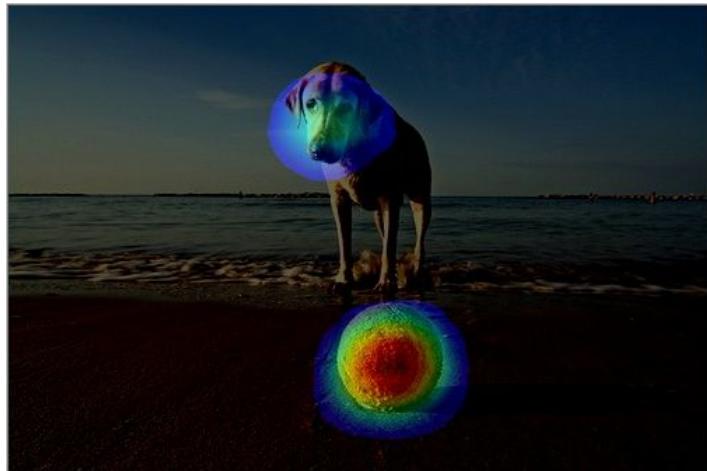
Obtaining Discriminative Regions



B. Zhou, A. Khosla, Lapedriza. A., A. Oliva, and A. Torralba. 2016. [Learning Deep Features for Discriminative Localization](#). CVPR (2016).

Class Activation Maps

Examples



Tennis Ball



Weimaraner

Simple Classes

Class Activation Maps

Examples



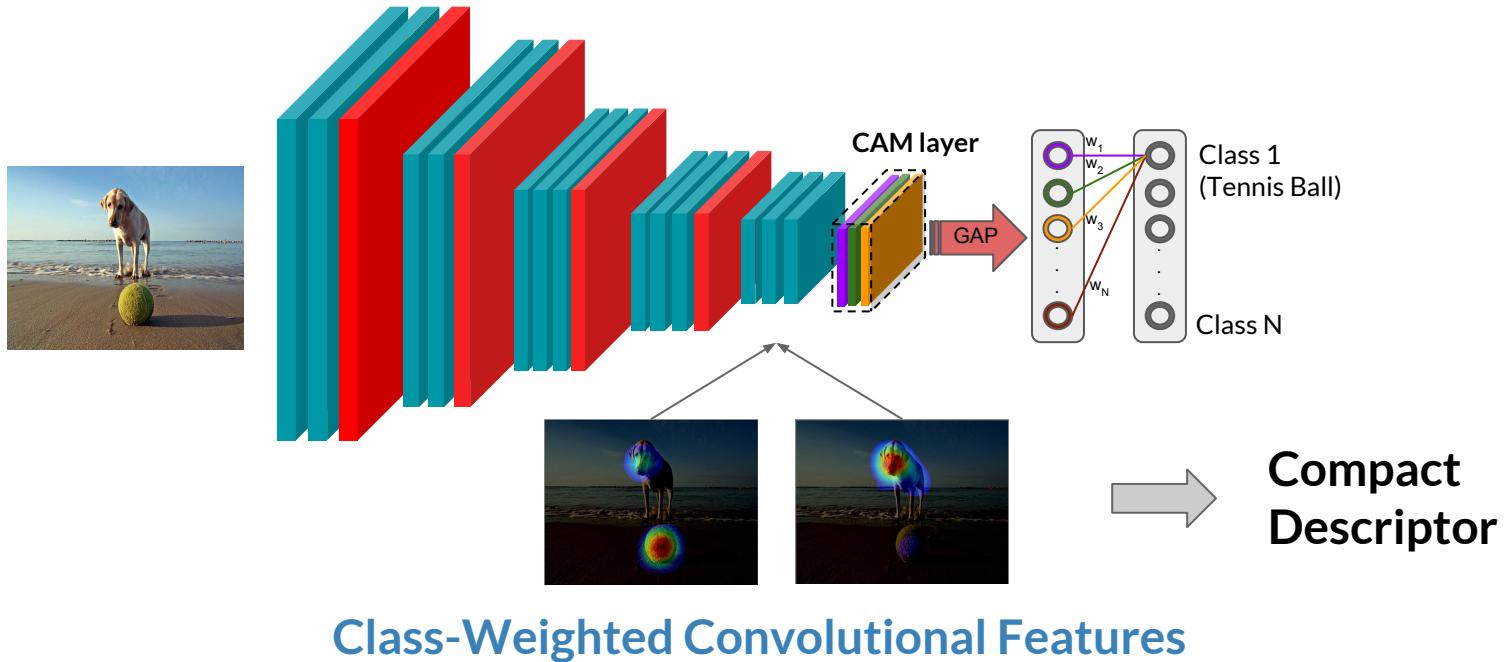
Sand Bar



Seashore

Complex Classes

Proposal



Class-Weighted Convolutional Features

Encode images dynamically combining **their semantics** with **convolutional features** using only the **knowledge contained** inside a **trained network**

Image Encoding Pipeline

1. Features & CAMs Extraction

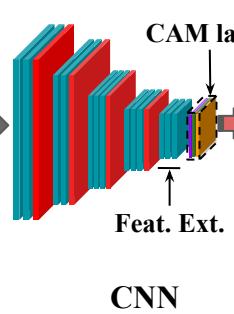
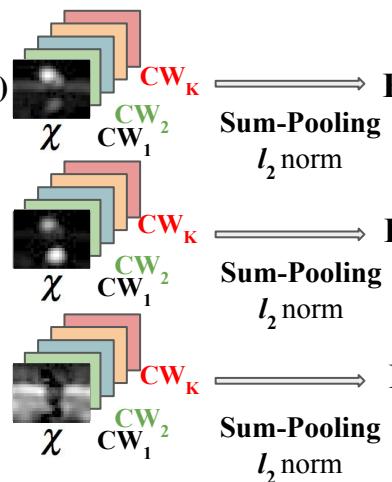


Image: I

2. Feature Weighting and Pooling



3. Descriptor Aggregation

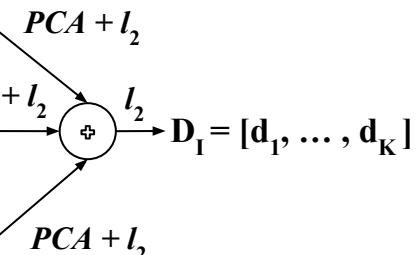
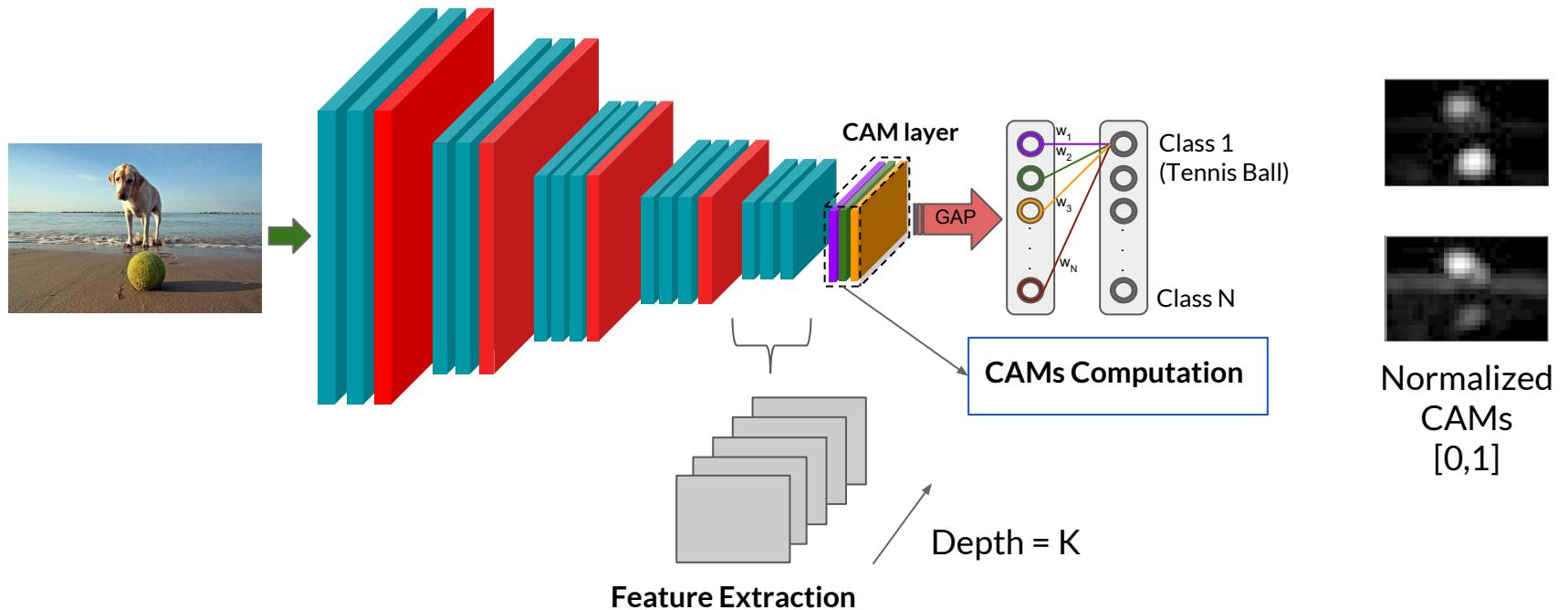


Image Encoding Pipeline

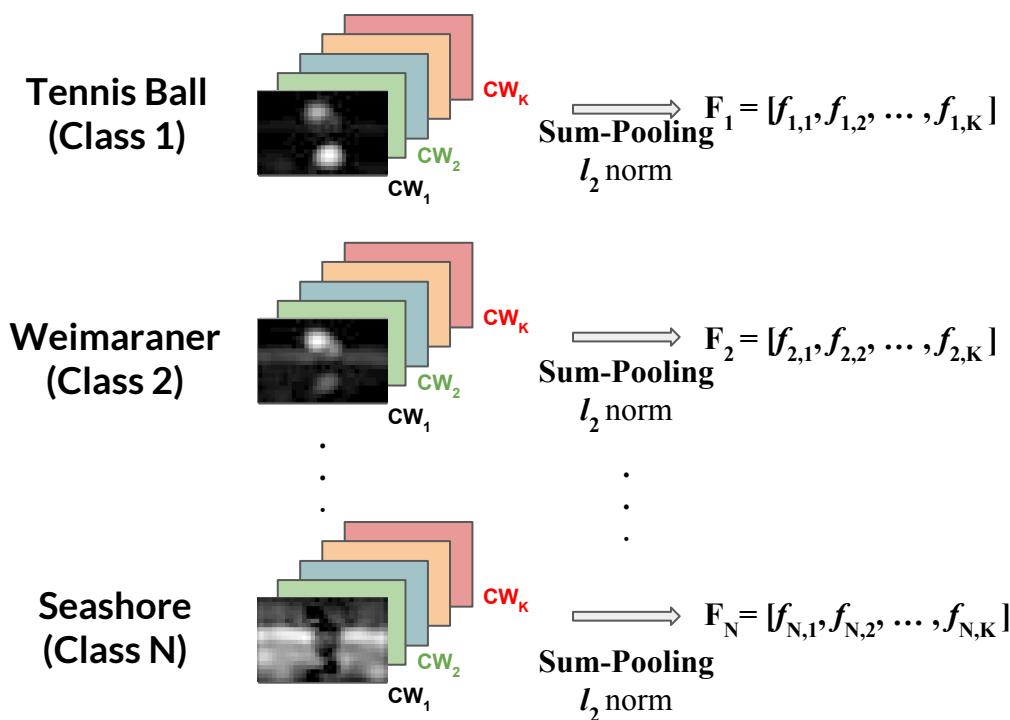
1. Features & CAMs Extraction



In a **single forward pass** we extract convolutional features and image CAMs

Image Encoding Pipeline

2. Feature Weighting and Pooling



One vector per class:

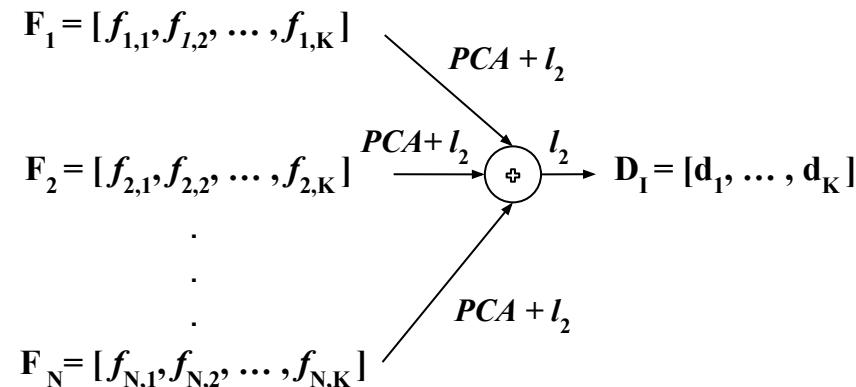
$$F^c = [f_1^c, f_2^c, \dots, f_K^c]$$

$$f_k^{(c)} = CW_k \sum_{i=1}^W \sum_{j=1}^H \chi_{i,j}^{(k)} CAM_{i,j}^{(c)}$$

Spatial Weighting using Class Activation Maps
Channel Weighting based on feature maps sparsity (as in CroW)

Image Encoding Pipeline

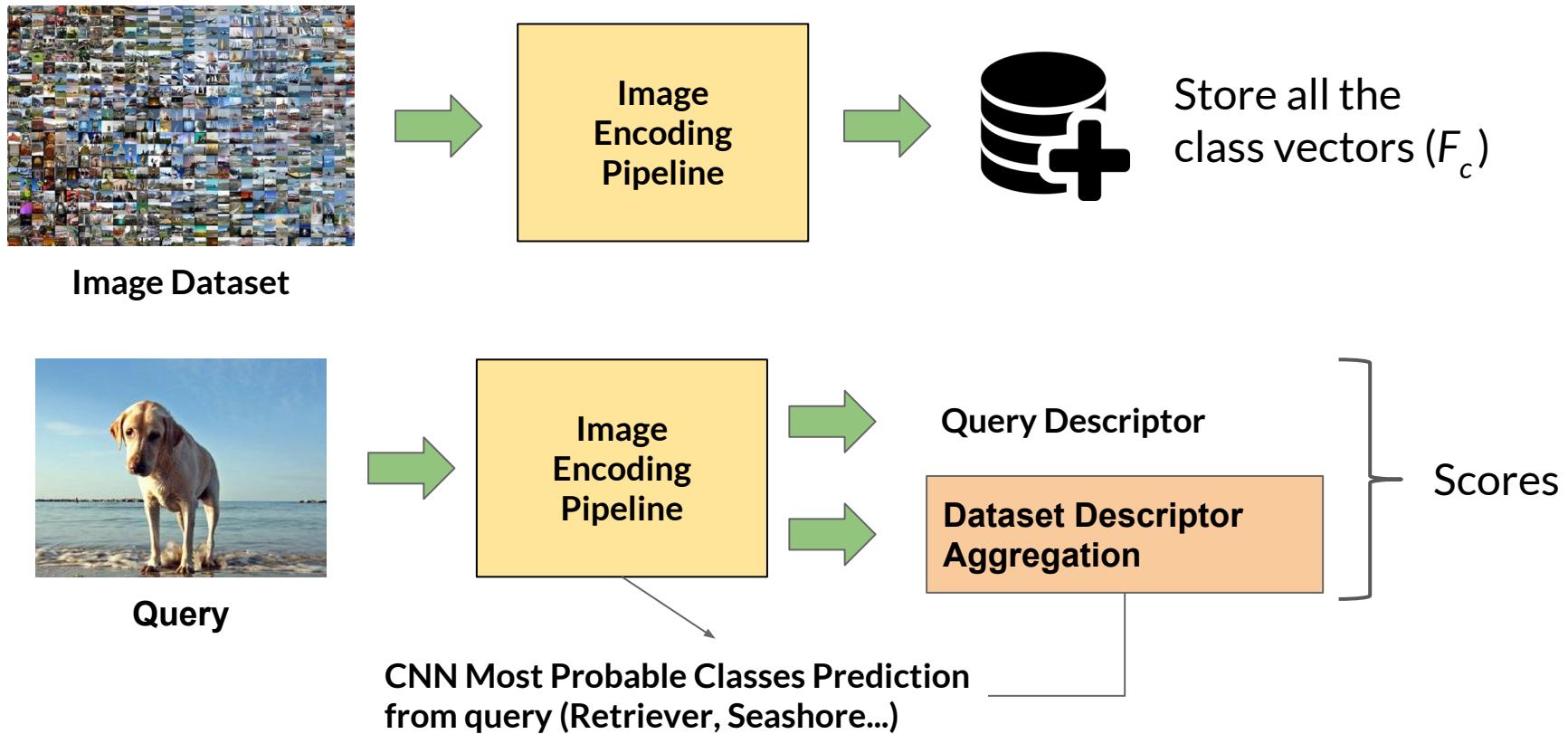
3. Descriptor Aggregation



- Which classes do we aggregate?
- What is the optimal number of classes (N_C)?
- What is the optimal number of classes (N_{PCA}) to compute the PCA ?

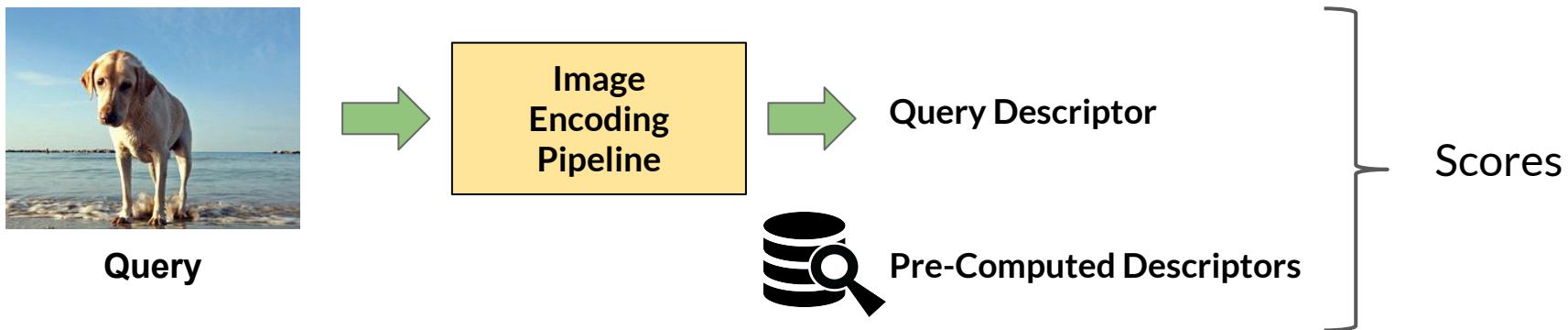
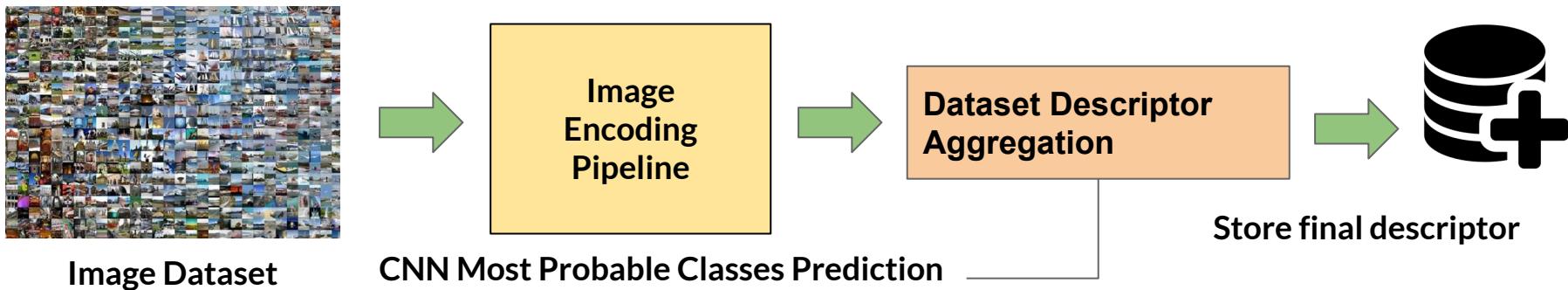
Descriptor Aggregation Strategies

Online Aggregation (OnA)



Descriptor Aggregation Strategies

Offline Aggregation (OfA)



Experiments

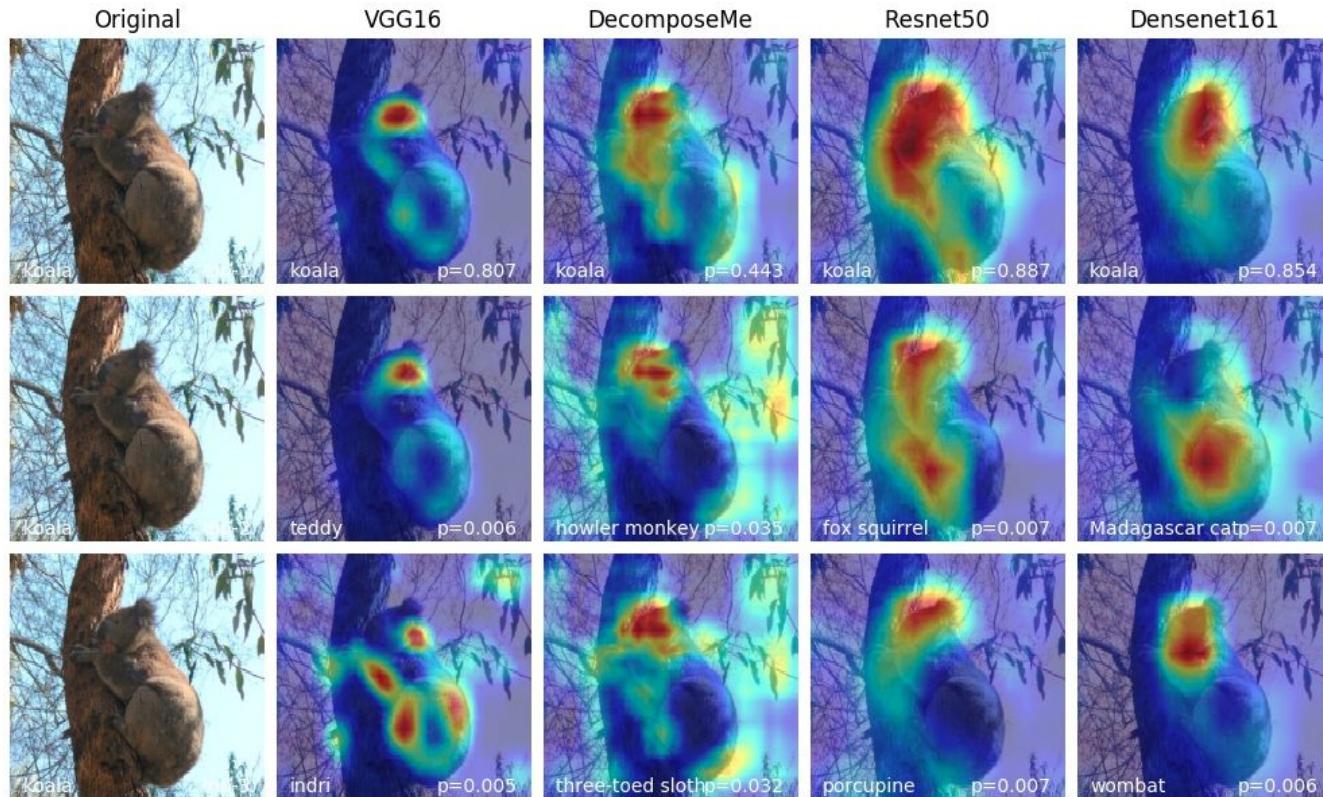
Experimental Setup

- Framework: Keras with Theano as backend / PyTorch
- Images resized to 1024x720 (keeping aspect ratio)
- We explored using different CNNs as feature extractors



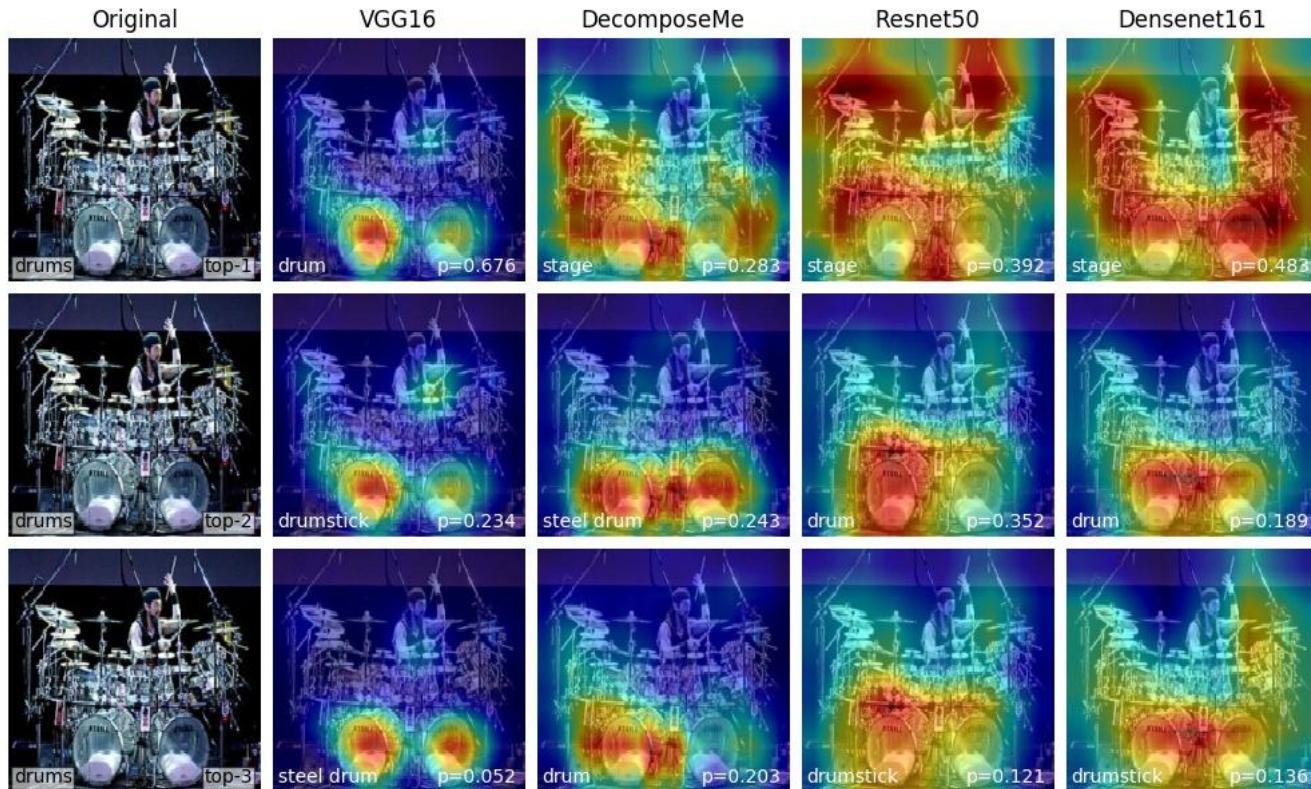
Experimental Setup

Examples of CAMs from different networks



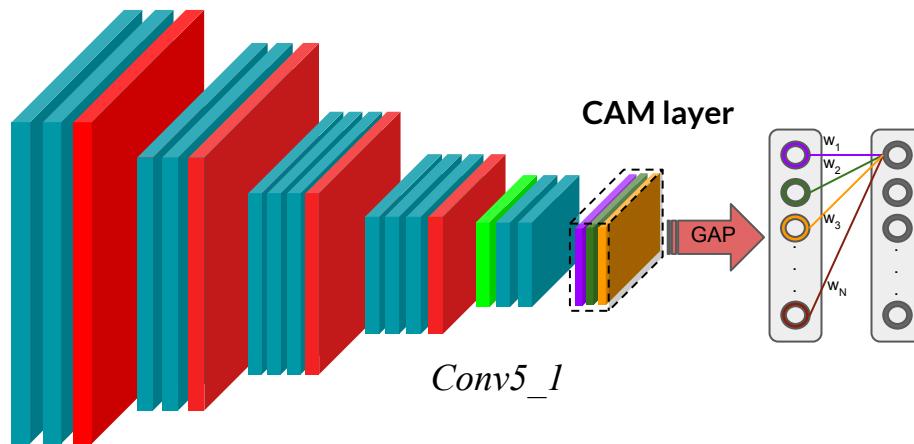
Experimental Setup

Examples of CAMs from different networks



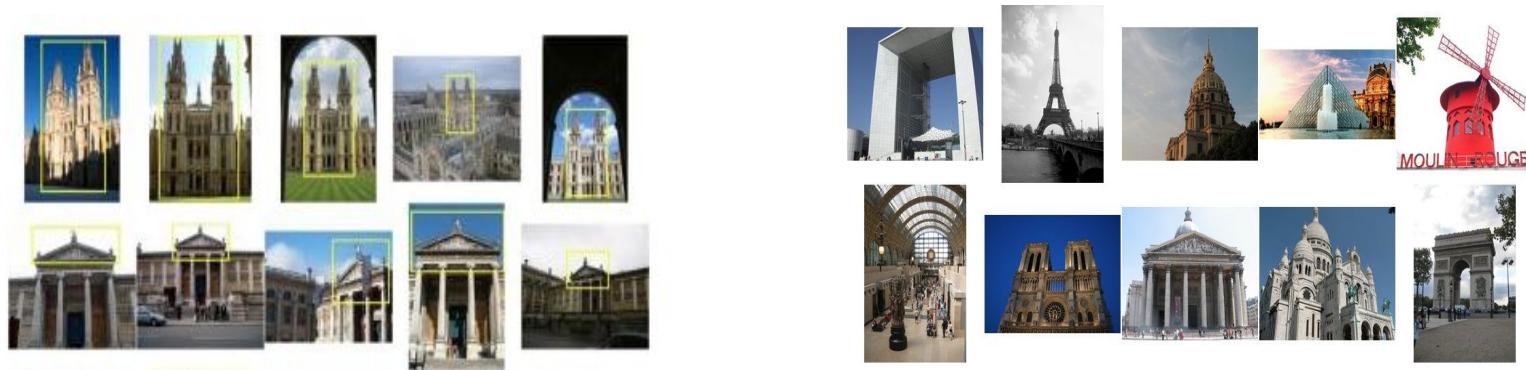
Experimental Setup

- VGG-16 CAM model as feature extractor (pre-trained on ImageNet)
- Features from *Conv5_1* Layer



Experimental Setup

- Datasets
 - Oxford 5k
 - Paris 6k
 - 100k Distractors (Flickr) → Oxford105k, Paris106k
- Using the features that fall inside the cropped query
- PCA computed with Oxf5k when testing in Par6k and vice versa



Philbin, J., Chum, O., Isard, M., Sivic, J. and Zisserman, A. [Object retrieval with large vocabularies and fast spatial matching](#), CVPR 2007
Philbin, J., Chum, O., Isard, M., Sivic, J. and Zisserman, A. [Lost in Quantization: Improving Particular Object Retrieval in Large Scale Image Databases](#). CVPR 2008

Experimental Setup

- Scores computed with cosine similarity

$$(\|X\|_2 =, \|Y\|_2 = 1)$$

$$\text{Cosine similarity } (X, Y) = \frac{X \cdot Y}{\|X\|_2 \|Y\|_2} = X \cdot Y$$

- Evaluation metric: mean Average Precision (mAP)

Ablation Studies

Finding the optimal parameters

- We carried out ablation studies to validate our method
- Computational Burden
- Number of top-classes to add: N_c
- Number of top-classes to compute the PCA matrix: N_{PCA}

Ablation Studies

Baseline Results & Computational Burden

Descriptor Aggregation	Oxford5k	Paris6k
Raw Features	0.396	0.526
Raw + Crow (channel)	0.420	0.549
Raw Features + PCA	0.589	0.662
Raw + Crow(channel) + PCA	0.607	0.685

Descriptor Aggregation	Time (s)	mAP
Raw + PCA	0.49	0.589
1 CAM	0.5	0.667
8 CAMs	0.6	0.709
32 CAMs	0.9	0.711
64 CAMs	1.5	0.712

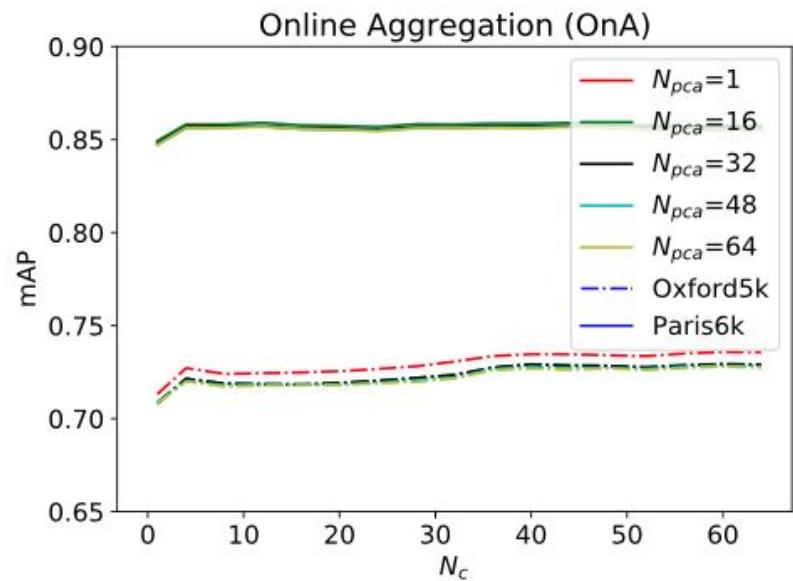
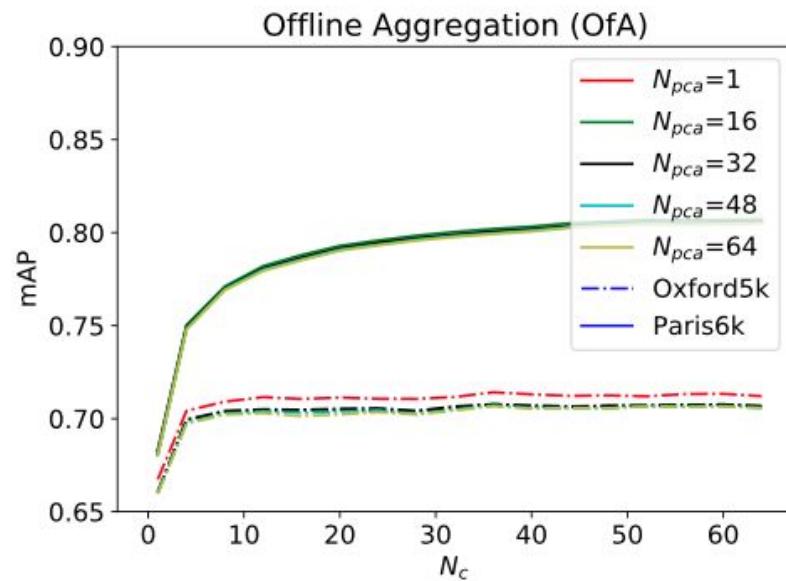
Ablation Studies

Networks Comparison

Network	Oxf5k	Paris6k
VGG-16 (Raw)	0.396	0.526
VGG-16 (64CAMs)	0.712	0.805
Resnet-50 (Raw)	0.389	0.508
Resnet-50 (64CAMs)	0.699	0.804
Densenet-161 (Raw)	0.339	0.495
Densenet-161 (64CAMs)	0.695	0.799

Ablation Studies

Sensitivity with respect to N_c and N_{pca}

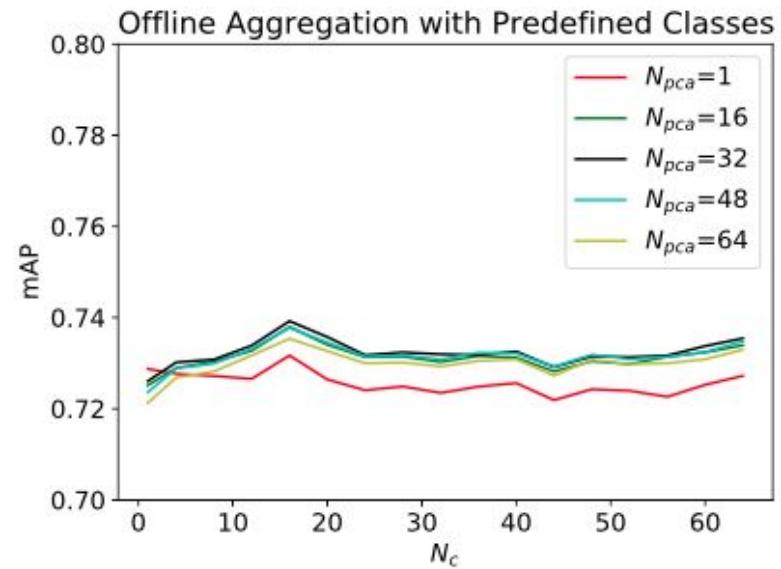
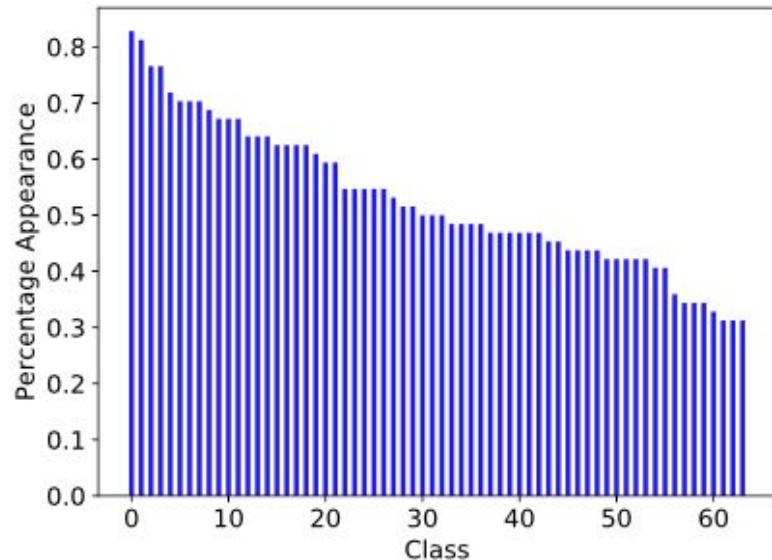


N_c corresponds to the number of classes aggregated.

N_{pca} corresponds to the classes used to compute the PCA (computed with data of Oxford when testing in Paris and vice versa).

Ablation Studies

Using predefined classes



The first 16 classes correspond to: *vault, bell cote, bannister, analog clock, movie theater, coil, pier, dome, pedestal, flagpole, church, chime, suspension bridge, birdhouse, sundial, triumphal arch*. All of them are related to **buildings**.

Comparison with State-of-the-Art

Method	Dim	Oxf5k	Par6k	Oxf105k	Par106k
SPoC[3]	256	0.531	-	0.501	-
uCroW[14]	256	0.666	0.767	0.629	0,695
CroW[14]	512	0.682	0.796	0.632	0.710
R-MAC[31]	512	0.669	0.830	0.616	0.757
BoW[16]	25k	0.738	0.820	0.593	0.648
Razavian [22]	32k	0.843	0.853	-	-
Ours(OnA)	512	0.736	0.855	-	-
Ours(OfA)	512	0.712	0.805	0.672	0.733

$$N_c = 64, N_{\text{pca}} = 1$$

Search Post-Processing

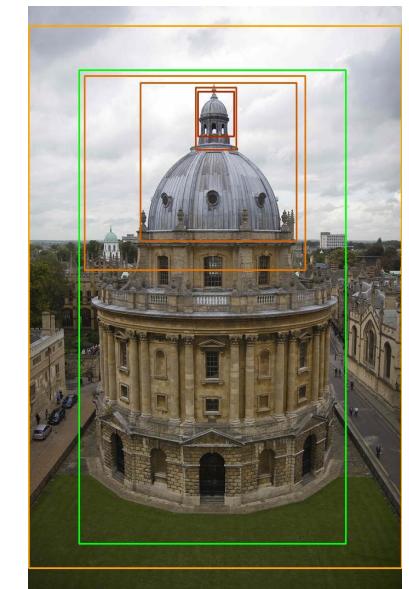
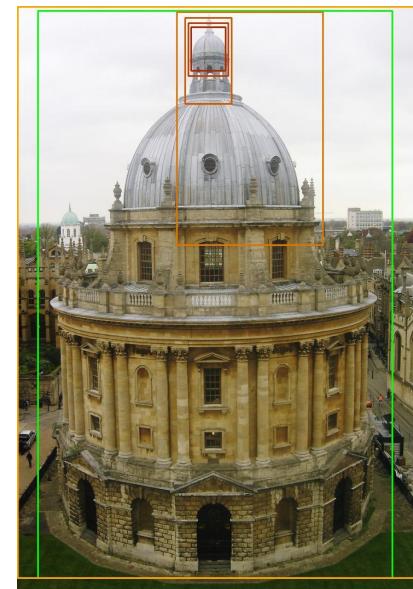
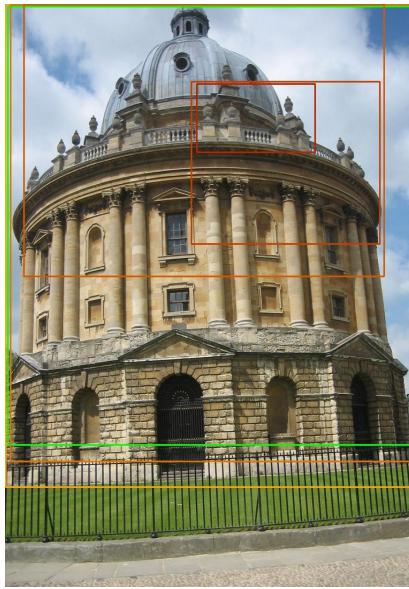
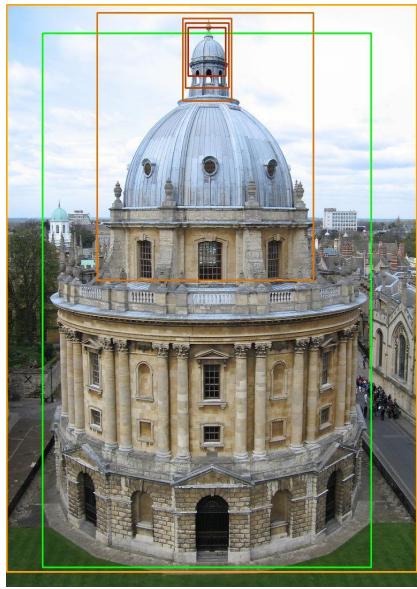
Region Proposals for Re-Ranking

- We establish a set of thresholds based on the max value of the normalized CAM
 - 1%, 10%, 20%, 30%, 40%
- Generate a bounding box covering the largest connected element
- Compare the descriptor of each new region of the target image with the query descriptor and generate new scores.
- Results a bit better if we average over 2 most probable CAMs
- For the top-R target images

Query Expansion

- Generate a new query by l_2 normalizing the sum of the top-QE target descriptors

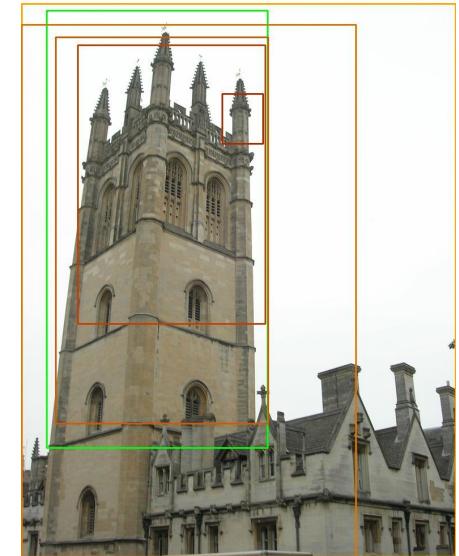
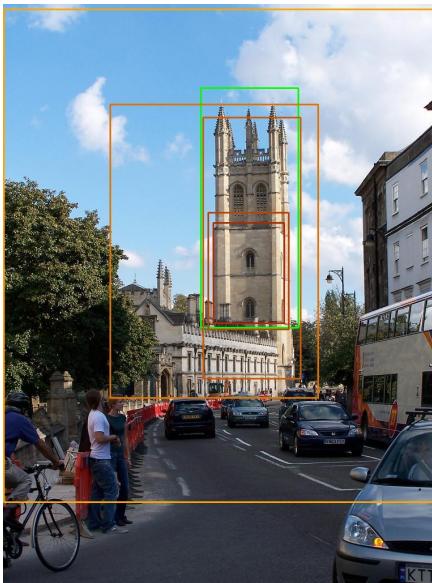
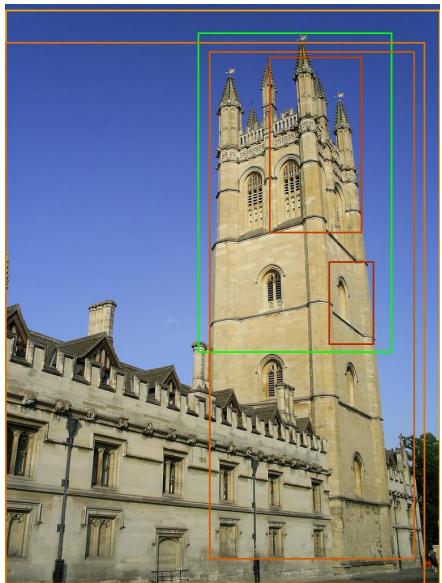
Region Proposals using CAMs



Green - Ground Truth

Orange → Red - Different Thresholds

Region Proposals using CAMs



Green - Ground Truth

Orange → Red - Different Thresholds

Comparison with State-of-the-Art

After Re-Ranking and Query Expansion

Method	Dim	R	QE	Oxf5k	Par6k	Oxf105k	Par106k
CroW	512	-	10	0.722	0.855	0.678	0.797
Ours(OnA)	512	-	10	0.760	0.873	-	-
Ours(OfA)	512	-	10	0.730	0.836	0.712	0.791
BoW	25k	100	10	0.788	0.848	0.651	0.641
Ours(OnA)	512	100	10	0.780	0.874	-	-
Ours(OfA)	512	100	10	0.773	0.838	0.750	0.780
RMAC	512	1000	5	0.770	0.877	0.726	0.817
Ours(OnA)	512	1000	5	0.811	0.874	-	-
Ours(OfA)	512	1000	5	0.801	0.855	0.769	0.800

$$N_c = 64, N_{pca} = 1, N_{re-ranking} = 6$$

Conclusions

Conclusions

- ▷ In this work we proposed a technique to build compact image representations focusing on their semantic content.
- ▷ We employed an image encoding pipeline that makes use of a pre-trained CNN and Class Activation Maps to extract discriminative regions from the image and weight its convolutional features accordingly
- ▷ Our experiments demonstrated that selecting the relevant content of an image to build the image descriptor is beneficial, and contributes to increase the retrieval performance. The proposed approach establishes a new state-of-the-art compared to methods that build image representations combining off-the-shelf features using random or fixed grid regions.

Thank
you

