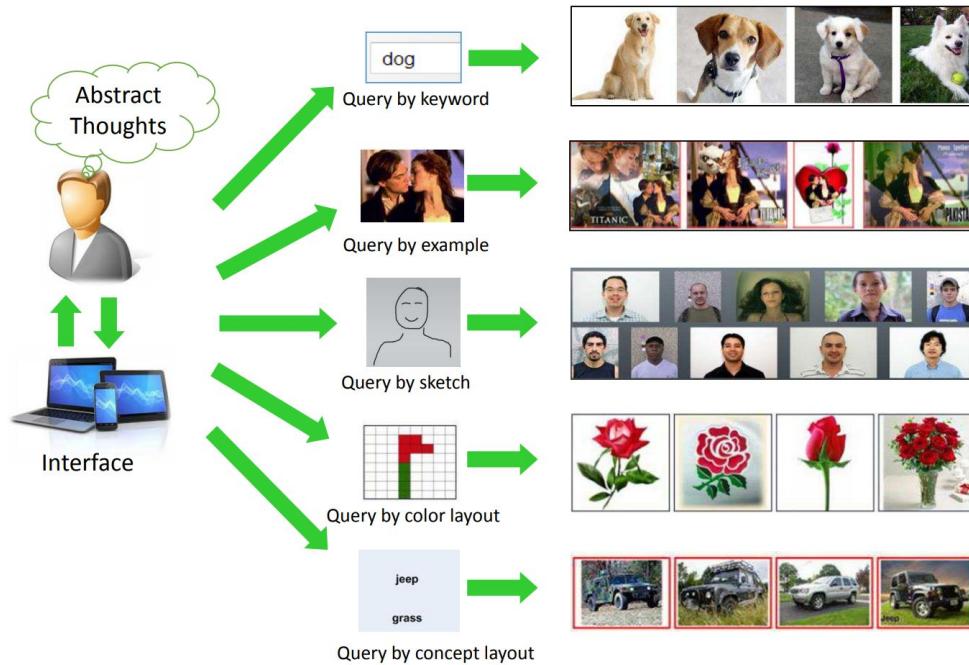


# Image Search

# Start of CBIR

5.6 Practical applications of CBIR.....	31
5.6.1 Crime prevention.....	31
5.6.2 The military .....	32
5.6.3 Intellectual property.....	32
5.6.4 Architectural and engineering design.....	33
5.6.5 Fashion and interior design.....	33
5.6.6 Journalism and advertising.....	33
5.6.7 Medical diagnosis .....	34
5.6.8 Geographical information systems (GIS) and remote sensing.....	34
5.6.9 Cultural heritage .....	35
5.6.10 Education and training.....	35
5.6.11 Home entertainment.....	35
5.6.12 Web searching.....	36
5.6.13 Conclusions.....	36
5.7 Current research trends .....	36

# survey CBIR



# Visual Geolocalization

# GeoWarp

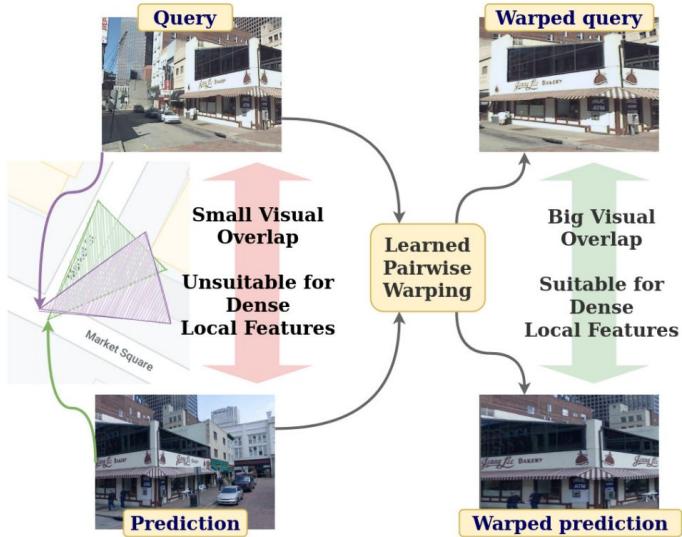
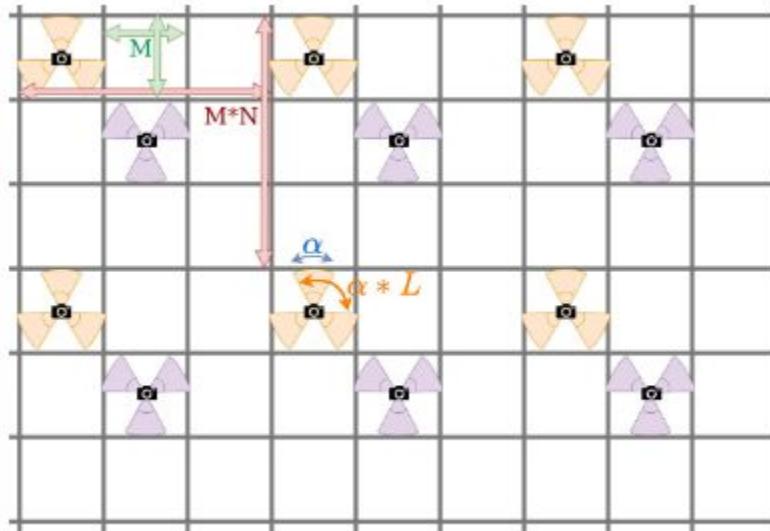
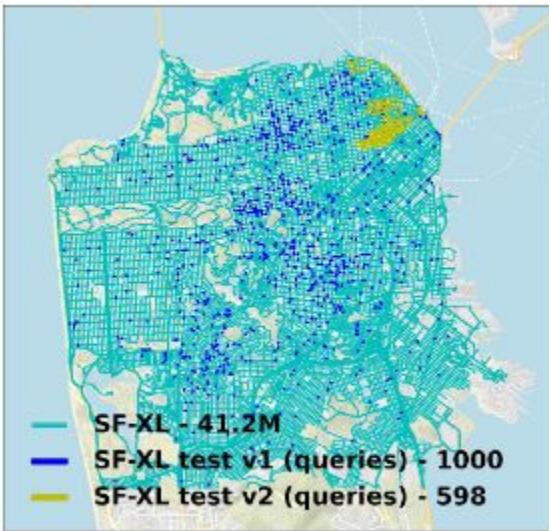


Figure 1. The appearance of two different views of the same place may differ significantly, thus making it hard to match them. Our method warps both images to a closer geometrical space and then computes their similarity using deep dense local features.

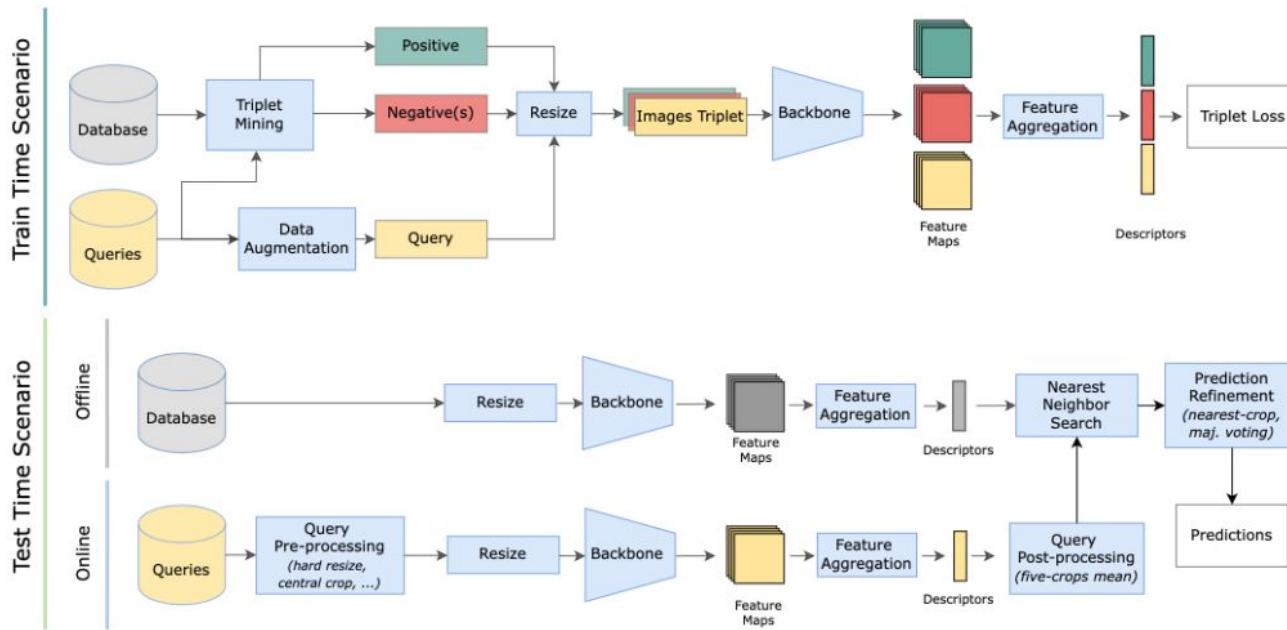
# Dataset: SF-XL



# St Lucia

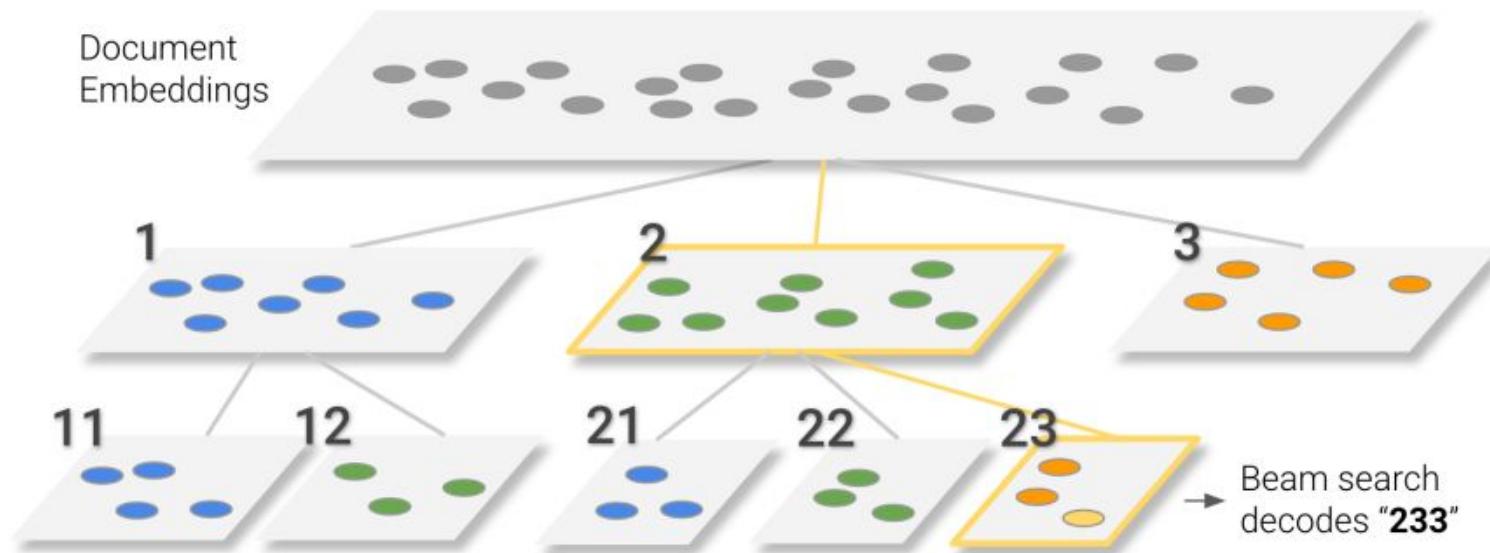
[Milford, et al. Mapping a Suburb With a Single Camera Using a Biologically Inspired SLAM System. TRO 2008](#)

# Benchmark



# Indexing

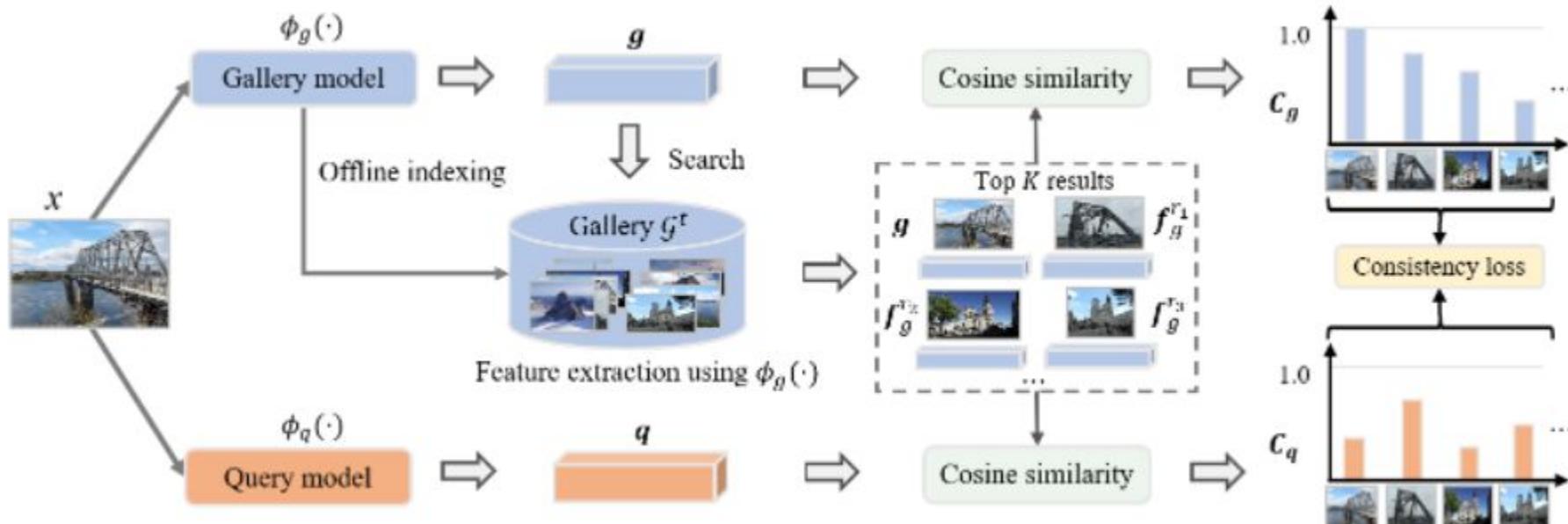
# Transformer Memory as a Differentiable Search Index



# Rerank

1

# Contextual Similarity



# **Total Recall: Automatic Query Expansion with a Generative Feature Model for Object Retrieval**

Ondřej Chum<sup>1</sup>, James Philbin<sup>1</sup>, Josef Sivic<sup>1</sup>, Michael Isard<sup>2</sup> and Andrew Zisserman<sup>1</sup>

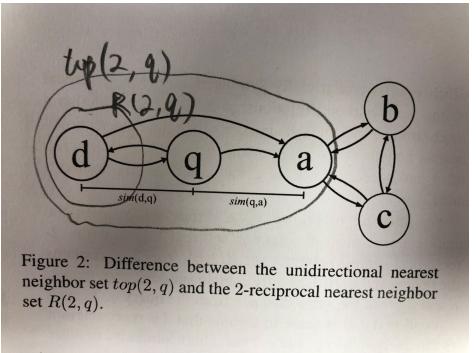
<sup>1</sup>Visual Geometry Group, Department of Engineering Science, University of Oxford

<sup>2</sup>Microsoft Research, Silicon Valley

{ondra, james, josef, az}@robots.ox.ac.uk

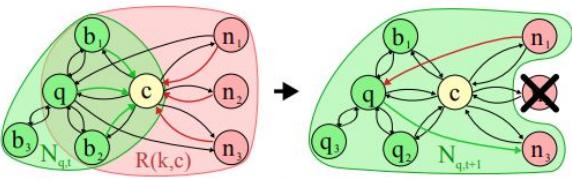
misard@microsoft.com

# K-reciprocal nearest neighbors

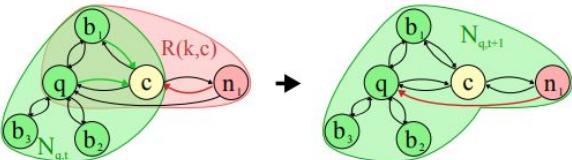


Make a close set for every query or every db image.

- 互相离得近的是 close set, 否则是 far set
- Far set中和close set近, 和其他图片远的  
比较重要



(a) Node C's neighborhood is considered, if it contains more than the half of the initial set, i.e.  $|N_{q,t} \cap R(k, c)| > \frac{1}{2} |N_{q,t}|$



(b) Node C's neighborhood is considered, if it adds less unknown nodes than known, i.e.  $|N_{q,t} \cap R(k, c)| > |R(k, c) \setminus N_{q,t}|$

---

## Algorithm 1: Expansion step

---

```

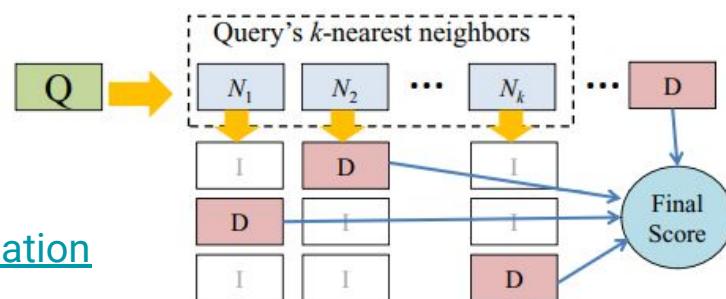
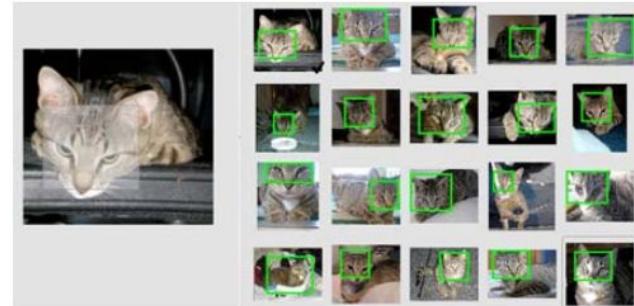
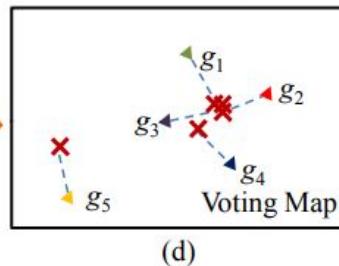
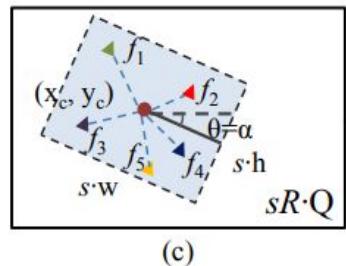
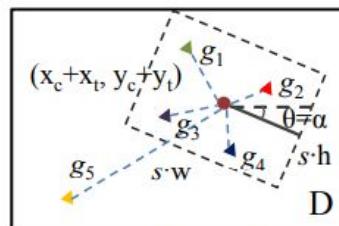
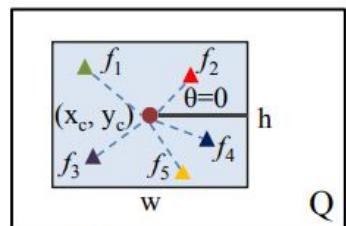
1 for  $t \leftarrow 0$  to 2 do
2    $N_{q,t+1} \leftarrow N_{q,t};$ 
3   foreach  $n \in N_{q,t}$  do
4     if  $|N_{q,t} \cap R(k, n)| > \frac{1}{2} |N_{q,t}|$  then
5        $N_{q,t+1} \leftarrow R(k, n) \cup N_{q,t+1};$ 
6     if  $|N_{q,t} \cap R(k, n)| > |R(k, n) \setminus N_{q,t}|$  then
7        $N_{q,t+1} \leftarrow R(k, n) \cup N_{q,t+1};$ 

```

---

The first condition allows only nodes which are connected to at least half of the *close set* to bring in their neighbors. This high connectivity ensures that added nodes are very likely to be relevant to the query. The second condition relaxes this restriction slightly by allowing weakly connected nodes to bring in their neighbors if the amount of new neighbors is smaller than the amount of connections already made to the *close set*. Nodes added to  $N_{q,t+1}$  are sorted according to  $sim(q, d)$  and inserted in this order into

# Spatially Constrained Similarity Measure

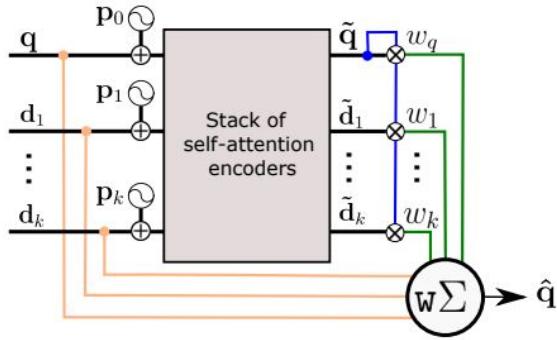


US8880563B2 Image search by query object segmentation

Shen, et al... Spatially-Constrained Similarity Measure for Large-Scale Object Retrieval, TPAMI 2014

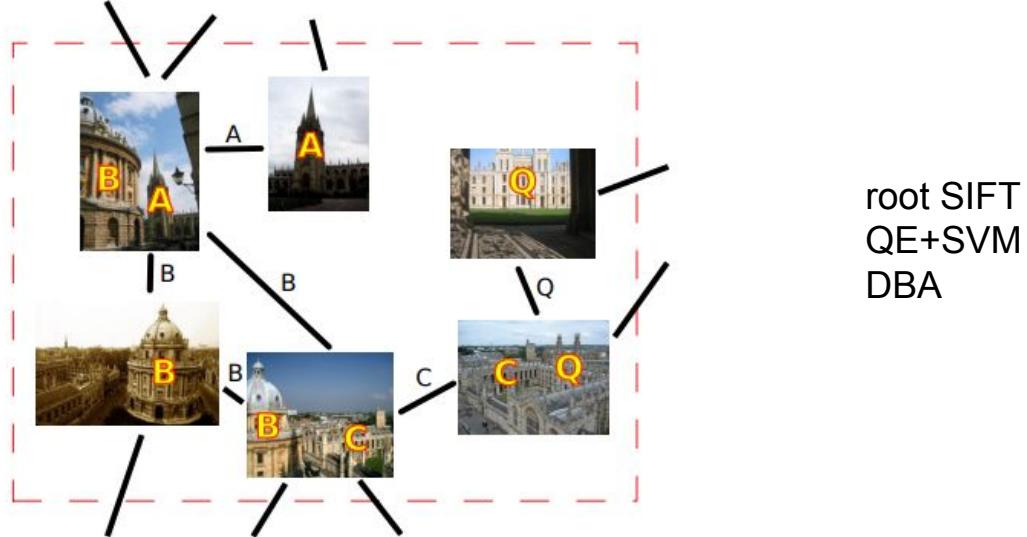
Shen, et al... Object retrieval and localization with spatially-constrained similarity measure and k-NN re-ranking, CVPR 2012 (210)

# Attention-Based Query Expansion

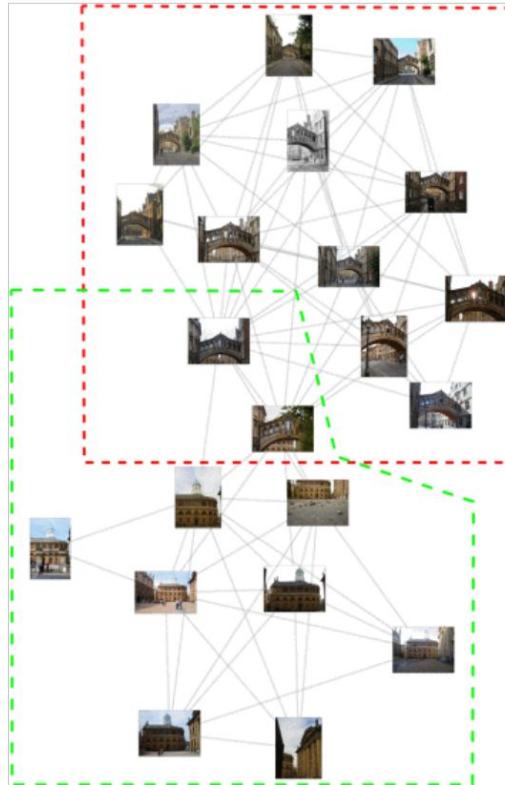


Method	$\theta(\mathbf{d}_i   \mathbf{q}, \{\mathbf{d}\}^k, \{\mathbf{d}\}^-, i) = w_i \mathbf{d}_i$
[7] <b>AQE:</b> Average QE	$w_i = 1$
[13] <b>AQEwD:</b> AQE with decay	$w_i = (k - i)/k$
[3] <b>DQE:</b> Discriminative QE	$\mathbf{w}$ is the dual-form solution of an SVM optimization problem using $\{\mathbf{d}\}^k$ as positives and $\{\mathbf{d}\}^-$ as negatives
[33] <b><math>\alpha</math>QE:</b> $\alpha$ -weighted QE	$w_i = \text{sim}(\mathbf{q}, \mathbf{d}_i)^\alpha$ , with $\alpha$ being a hyperparameter.

# Database Augmentation (before CNN)



# Matching graph



# QE

Local features: [Aggregate13Tolias],[Three12Arandjelovic]

# Geometric verification

3

# Improve Spatial verification



# Spatial verification

**Table 6.** (a) The three affine sub-groups compared in the spatial re-ranking. (b) Computing  $H$  as  $H_2^{-1}H_1$ , preserving "upness" for the 5 dof case.

Transformation	dof	Matrix
translation + isotropic scale	3	$\begin{bmatrix} a & 0 & t_x \\ 0 & a & t_y \end{bmatrix}$
translation + anisotropic scale	4	$\begin{bmatrix} a & 0 & t_x \\ 0 & b & t_y \end{bmatrix}$
translation + vertical shear	5	$\begin{bmatrix} a & 0 & t_x \\ b & c & t_y \end{bmatrix}$

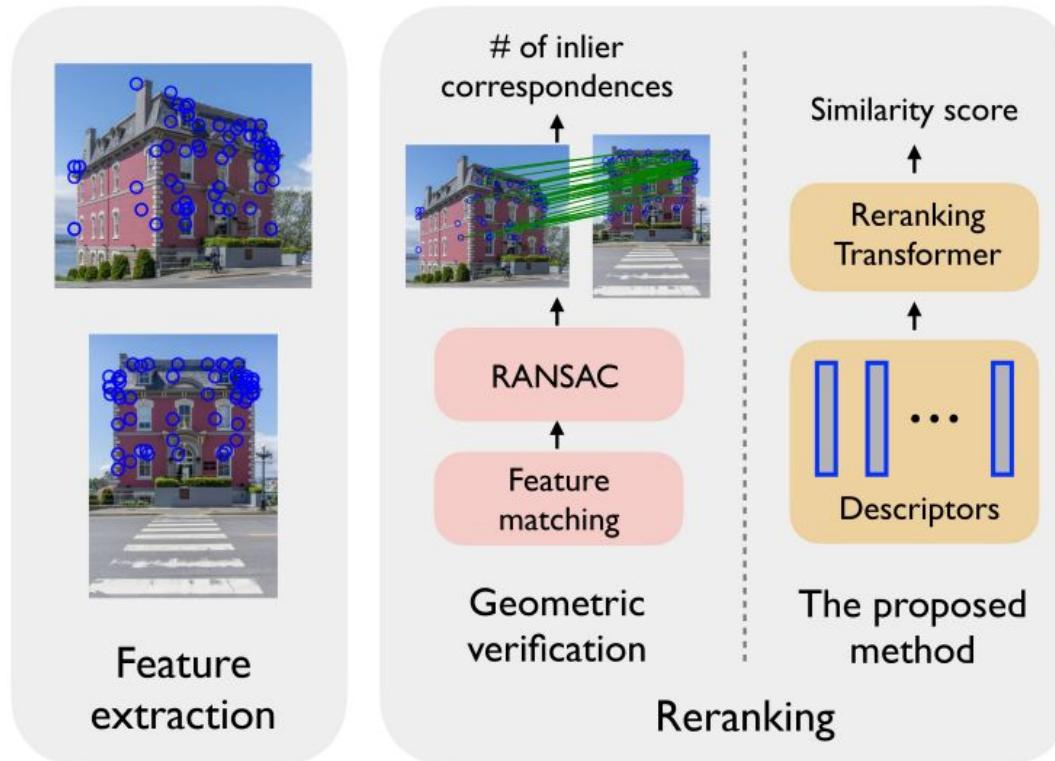
(a)

Diagram (a) illustrates three affine transformations based on their degrees of freedom (dof). The first row shows a transformation with 3 dof (translation + isotropic scale), represented by a matrix with entries  $a$  and  $t_x, t_y$ . The second row shows a transformation with 4 dof (translation + anisotropic scale), represented by a matrix with entries  $a, b$  and  $t_x, t_y$ . The third row shows a transformation with 5 dof (translation + vertical shear), represented by a matrix with entries  $a, b, c$  and  $t_x, t_y$ .

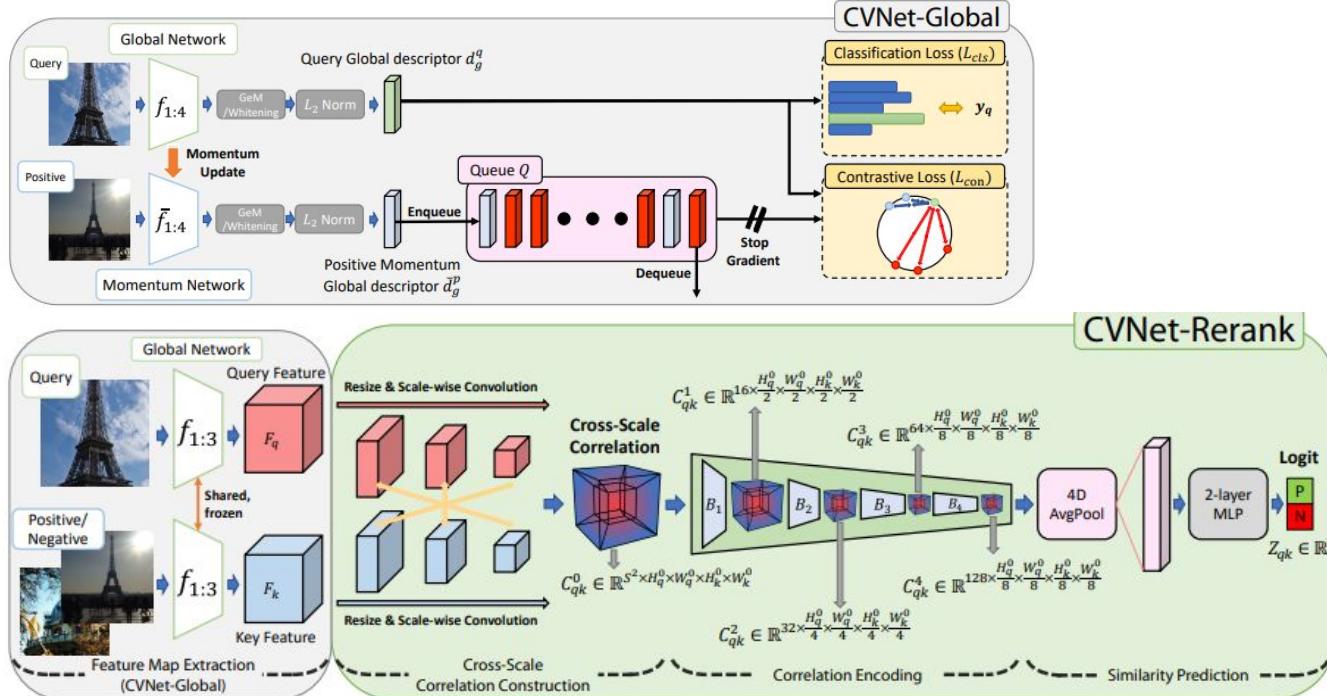
(b)

Diagram (b) illustrates the computation of the transformation  $H$  as  $H_2^{-1}H_1$  for the 5 dof case, while preserving "upness". It shows two ellipses,  $C_1$  and  $C_2$ , representing objects in the scene. Ellipse  $C_1$  is transformed by  $H$  into ellipse  $C_2$ . Ellipse  $C_1$  is also transformed by  $H_1$  into a circle. Ellipse  $C_2$  is also transformed by  $H_2$  into a circle. Finally, the circle resulting from  $H_1$  is transformed by  $I$  (the identity transformation) into the circle resulting from  $H_2$ . This diagram visually represents how the "upness" of features is preserved through the intermediate transformation  $H_1$ .

# Reranking transformers



# Correlation Verification





# Diffusion

3

# Diffusion in semi-supervised learning

$$F(t+1) = \alpha S F(t) + (1 - \alpha) Y$$

$$F^* = \lim_{t \rightarrow \infty} F(t) = (1 - \alpha)(I - \alpha S)^{-1} Y,$$

贡献是提出使  
用  $S = D^{-1/2} W D^{-1/2}$

Rigorous proof

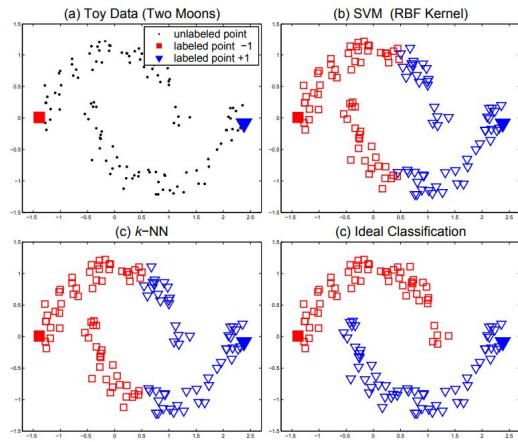


Figure 1: Classification on the two moons pattern. (a) toy data set with two labeled points; (b) classifying result given by the SVM with a RBF kernel; (c)  $k$ -NN with  $k = 1$ ; (d) ideal classification that we hope to obtain.

$$\begin{aligned}
 f(0) &= Y \\
 f(1) &= \alpha S Y + (1 - \alpha) Y \\
 f(2) &= (\alpha S)^2 [ \alpha S Y + (1 - \alpha) Y ] + (1 - \alpha) Y \\
 &= (\alpha S)^2 Y + (\alpha S)^2 (1 - \alpha) Y + (1 - \alpha) Y \\
 &= (\alpha S)^2 Y + (1 - \alpha) [ 1 + \alpha S ] Y \\
 f(3) &= \alpha S [ (\alpha S)^2 Y + (1 - \alpha) [ 1 + \alpha S ] Y ] + (1 - \alpha) Y \\
 &= (\alpha S)^3 Y + \alpha S (1 - \alpha) [ 1 + \alpha S ] Y + (1 - \alpha) Y \\
 &= (\alpha S)^3 Y + (1 - \alpha) [ \alpha S + \alpha S^2 ] Y + (1 - \alpha) Y \\
 &= (\alpha S)^3 Y + (1 - \alpha) [ 1 + \alpha S + \alpha S^2 ] Y
 \end{aligned}$$

# Diffusion in Image/Text Ranking

贡献是提出使

$$S = D^{-1/2} W D^{-1/2}$$

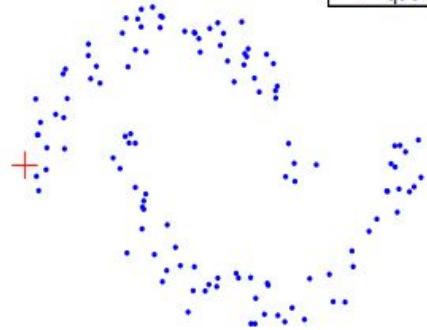
$$F(t+1) = \alpha S F(t) + (1 - \alpha) Y$$

$$F^* = \lim_{t \rightarrow \infty} F(t) = (1 - \alpha)(I - \alpha S)^{-1} Y,$$

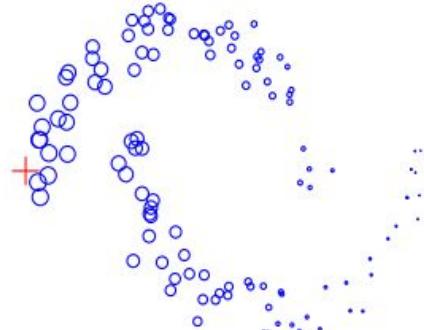
Easy proof

(a) Two moons ranking problem

+ query



(b) Ranking by Euclidean distance



(c) Ideal ranking

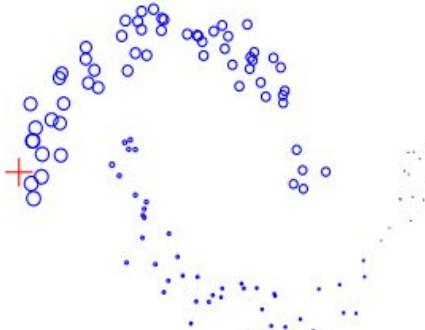


Figure 1: Ranking on the two moons pattern. The marker sizes are proportional to the ranking in the last two figures. (a) toy data set with a single query; (b) ranking by the Euclidean distances; (c) ideal ranking result we hope to obtain.

# Diffusion and Personalized PageRank

PageRank

$$P = (1 - \epsilon)U + \epsilon D^{-1}W,$$

Personalized PageRank

$$P = D^{-1}W.$$

$$\pi(t+1) = \alpha P^T \pi(t) + (1 - \alpha)y.$$

区别是  $U$  是 uniformed,  $y$  可以给感兴趣的网页赋予较大的值

Note that  $P = D^{-1}W$ .  $P = (1 - \epsilon)U + \epsilon D^{-1}W$ ,  
have same stationary distribution

[Zhou, et al...Ranking on Data Manifold, NIPS 2004 \(812\)](#)

Personalized PageRank

$$P = D^{-1}W \quad \pi(t+1) = \alpha P^T \pi(t) + (1 - \alpha)y$$

$$\pi^* = (1 - \alpha)(I - \alpha P^T)^{-1}y \rightarrow \text{见 Zhou 2003}$$

$$\pi^* = (I - \alpha P^T)^{-1}y \rightarrow \text{as } (1 - \alpha) \text{ has no influence}$$

$$= [(D - \alpha W)D^{-1}]^{-1}y$$

$$= \underbrace{D}_{\sim} (D - \alpha W)^{-1}y$$

Diffusion

$$S = D^{-\frac{1}{2}}WD^{-\frac{1}{2}} \quad f(t+1) = \alpha Sf(t) + (1 - \alpha)y$$

$$f^* = (1 - \alpha)(I - \alpha S)^{-1}y$$

$$f^* = (I - \alpha S)^{-1}y$$

$$= [I - \alpha D^{-\frac{1}{2}}WD^{-\frac{1}{2}}]^{-1}y$$

$$= [D^{-\frac{1}{2}}(D - \alpha W)D^{-\frac{1}{2}}]^{-1}y$$

$$= \underbrace{D^{\frac{1}{2}}}_{\sim} (D - \alpha W)^{-1} \underbrace{D^{\frac{1}{2}}}_{\sim} y$$

# Graph Transduction

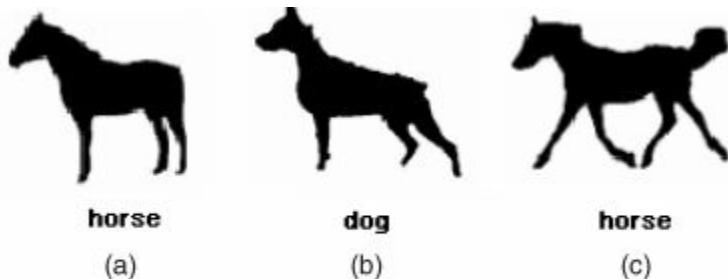


Fig. 1. Existing shape similarity methods incorrectly rank shape (b) as more similar to (a) than (c).

the original similarities  $w_{i,j} = \text{sim}(x_i, x_j)$ . Our intuition is that the new similarity  $f(x_i) = s(x_1, x_i)$  will be large iff all points  $x_j$  that are very similar to  $x_i$  (large  $\text{sim}(x_i, x_j)$ ) are also very similar to query  $x_1$  (large  $\text{sim}(x_1, x_j)$ ). Note that

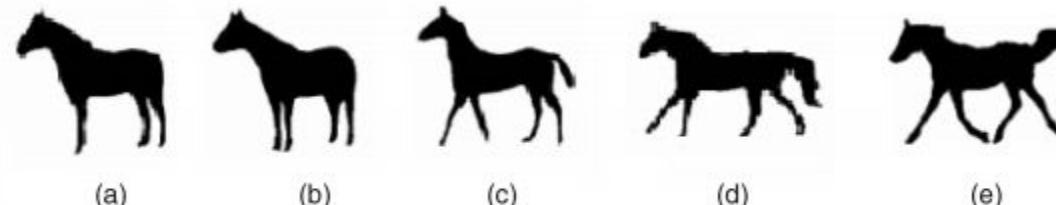


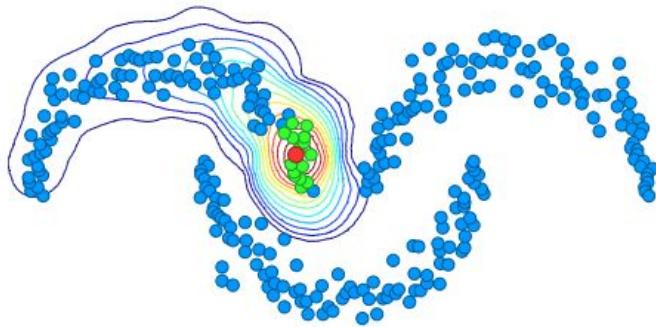
Fig. 2. A key idea of the proposed distance learning is to replace the original shape distance between (a) and (e) with a distance induced by geodesic paths in the manifold of known shapes. One such path is (a)-(e) in this figure.

# Summary of diffusion processes

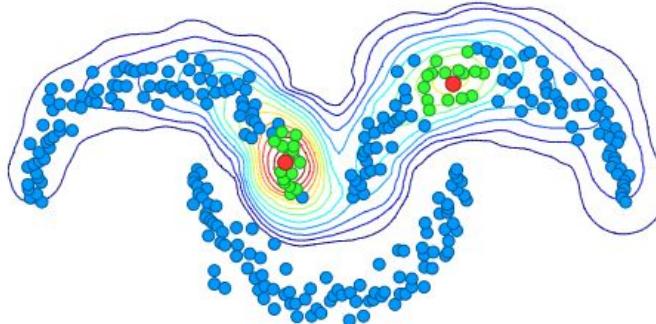
This paper uses clear math terms, so it is helpful to search and review the mathematical knowledge.

Method	Abbr.	Initialization $\mathbf{W}^0$	Transition $\mathbf{T}$	Diffusion
Global PageRank [17]	GPR	$\mathbf{u}$	$\mathbf{P}$	$\mathbf{f}_{t+1} = \mathbf{f}_t \mathbf{T}$
Personalized PageRank [17]	PPR	$\mathbf{u}$	$\mathbf{P}$	$\mathbf{f}_{t+1} = \alpha \mathbf{f}_t \mathbf{T} + (1 - \alpha) \mathbf{y}$
Ranking on Manifolds [23]	ROM	$\mathbf{u}$	$\mathbf{P}_{NC}$	$\mathbf{f}_{t+1} = \alpha \mathbf{f}_t \mathbf{T} + (1 - \alpha) \mathbf{y}$
Label Propagation [24]	LP	$\mathbf{y}$	$\mathbf{P}$	$\mathbf{f}_{t+1} = \mathbf{f}_t \mathbf{T}$ and $f(i) = 1$
Graph Transduction [2]	GT	$\mathbf{y}$	$\mathbf{P}$	$\mathbf{f}_{t+1} = \mathbf{f}_t \mathbf{T}$ and $f(i) = 1$
Locally Constrained DP [21]	LCDP	$\mathbf{A}$	$\mathbf{P}_{kNN}$	$\mathbf{W}_{t+1} = \mathbf{T} \mathbf{W}_t \mathbf{T}^T$
Tensor Graph Diffusion [22]	TGD	$\mathbf{A}$	$\mathbf{P}_{DS}$	$\mathbf{W}_{t+1} = \mathbf{T} \mathbf{W}_t \mathbf{T}^T + \mathbf{I}$
Shortest Path Propagation [20]	SPP	$\mathbf{y}$	$\mathbf{P}_{SP}$	$\mathbf{f}_{t+1} = \mathbf{f}_t \mathbf{T}$
Self Smoothing Operator [8]	SSO	$\mathbf{A}$	$\mathbf{P}$	$\mathbf{W}_{t+1} = \mathbf{W}_t \mathbf{T}$
Self Diffusion [19]	SD	$\mathbf{A}$	$\mathbf{P}$	$\mathbf{W}_{t+1} = \mathbf{W}_t \mathbf{T} + \mathbf{I}$
Replicator Dynamics [18]	RD	$\mathbf{u}$	$\mathbf{A}$	$\mathbf{f}_{t+1} = \mathbf{f}_t \odot \mathbf{T} \mathbf{f}_t$ and $\mathbf{f}_{t+1} = \mathbf{f}_{t+1} /  \mathbf{f}_{t+1} $
Power Iteration Clustering [11]	PIC	$\mathbf{s}$	$\mathbf{P}$	$\mathbf{f}_{t+1} = \mathbf{T} \mathbf{f}_t$ and $\mathbf{f}_{t+1} = \mathbf{f}_{t+1} /  \mathbf{f}_{t+1} $
Authority Shift Clustering [3]	ASC	$\mathbf{P}_{PPR}$	$\mathbf{P}_{PPR}$	$\mathbf{W}_{t+1} = \mathbf{W}_t \mathbf{T}$

# Efficient; region manifolds; small objects

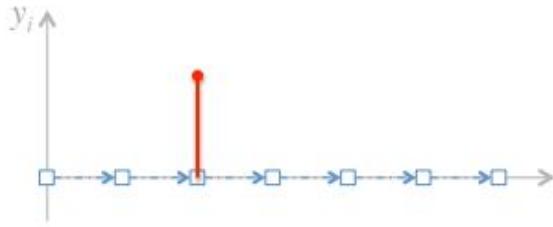


(a) single query

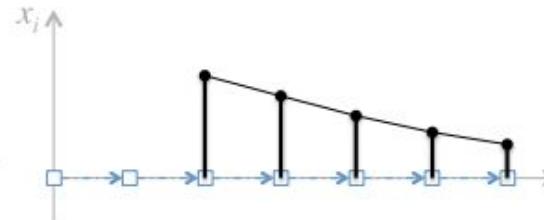


(b) multiple queries

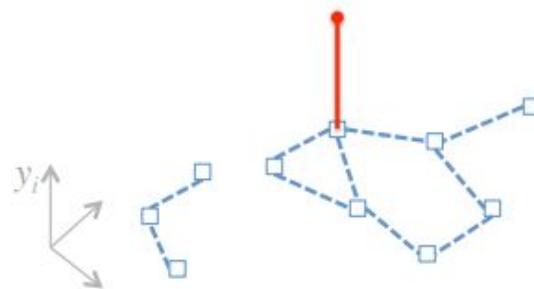
# Fast spectral ranking



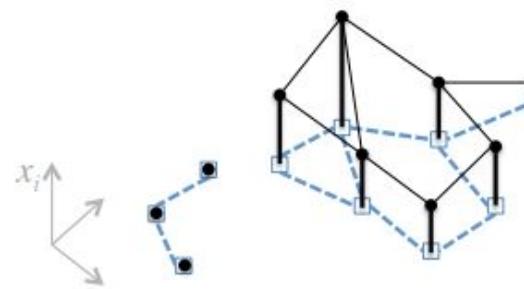
(a) Input signal  $y$



(b) Output signal  $x$

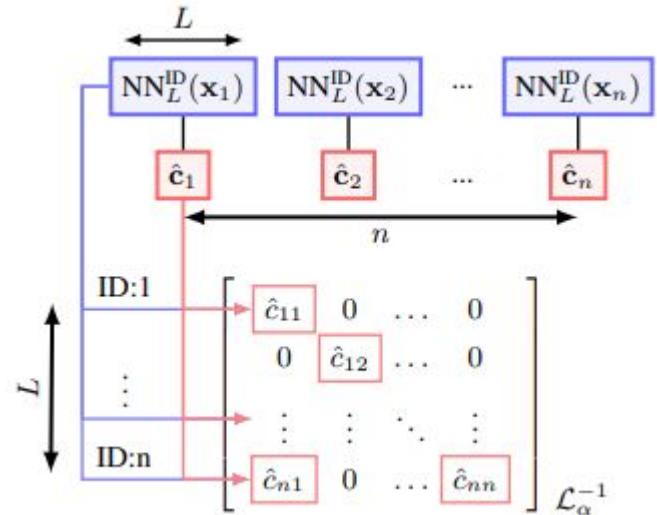


(c) Input signal  $y$

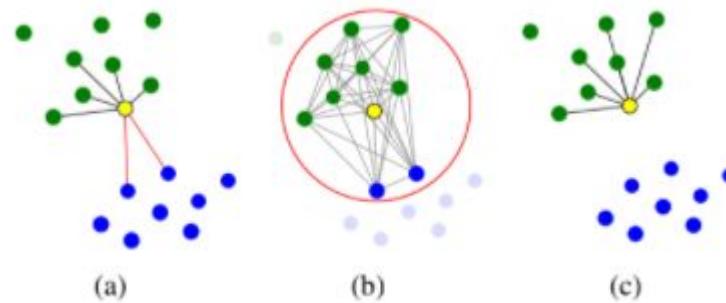
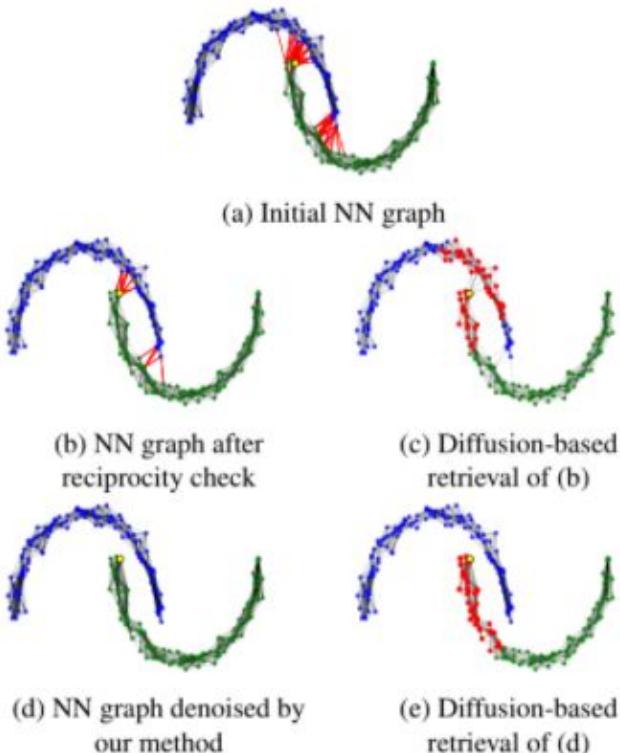


(d) Output signal  $x$

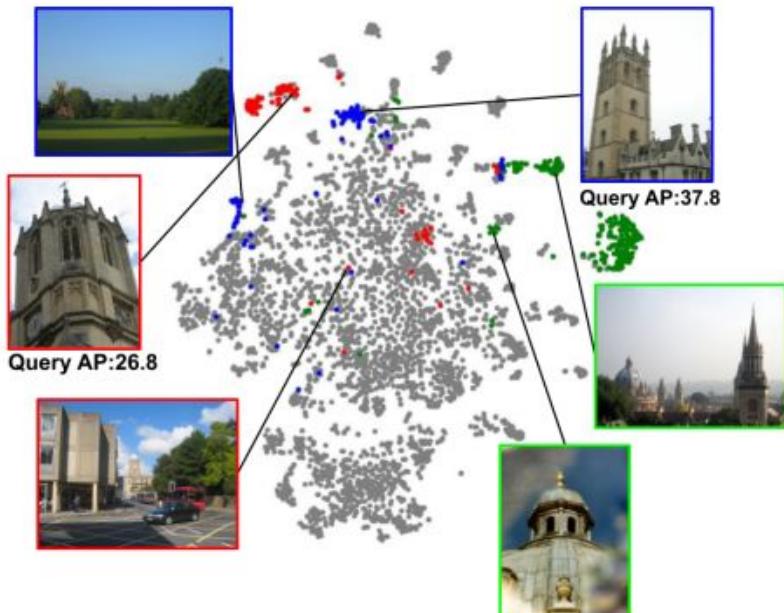
# Offline diffusion



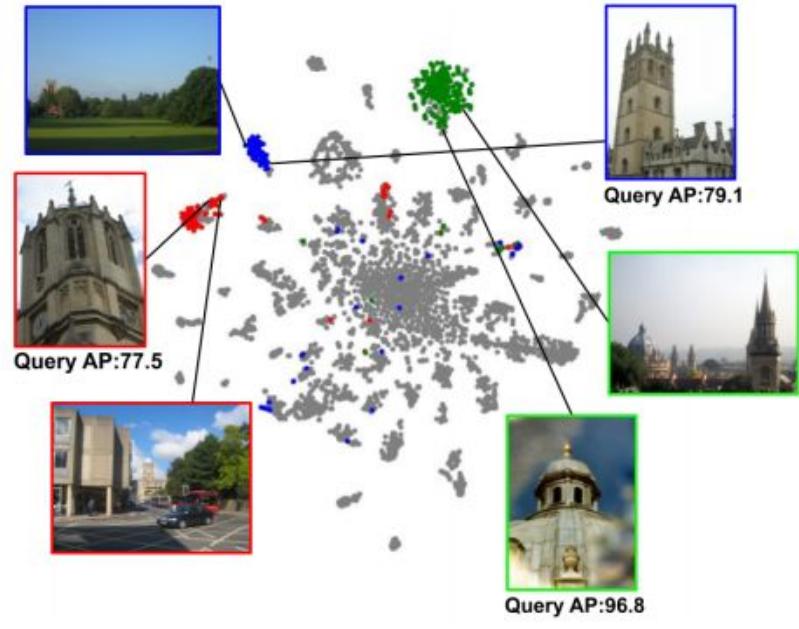
# Noisy edges



# Guided Similarity Separation



(a) GeM



(b) GeM+GSS

# Hypergraph diffusion

$$X(t) = X(0)P^t = X(0)AB\dots AB = X(0)AQ^{t-1}B$$

$$P = D_v^{-1}WD_e^{-1}W^T = AB.$$

: vertex-edge transition matrix  $A = D_v^{-1}W$   
edge-vertex transition matrix  $B = D_e^{-1}W^T$ .

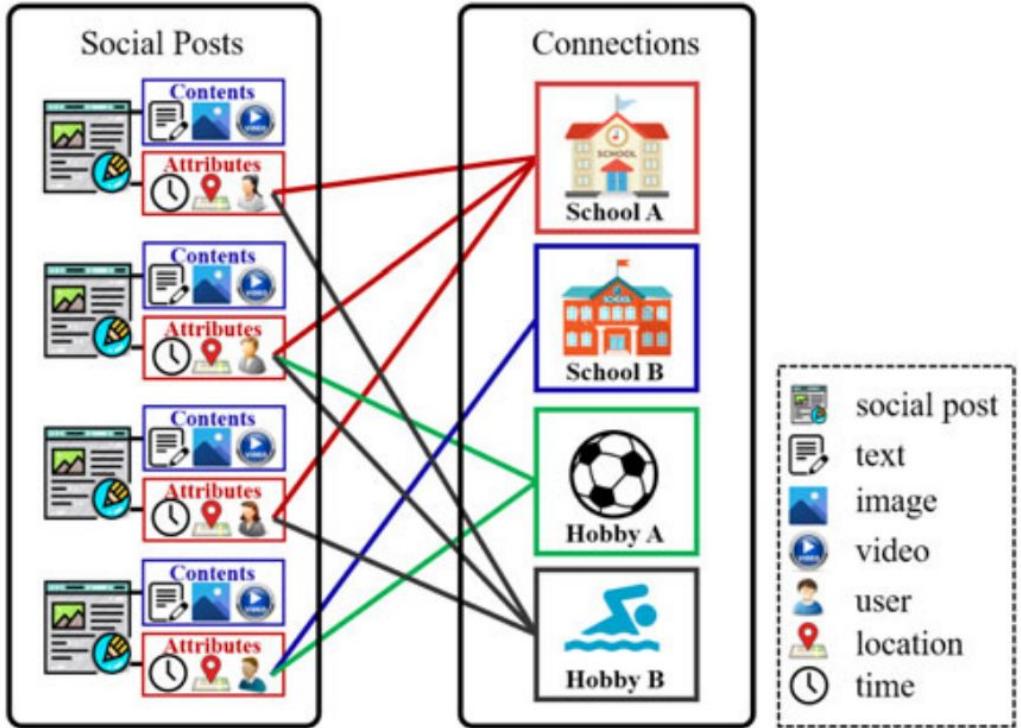
[Avin, et al...Radio cover time in hyper-graphs Ad Hoc Networks 2014](#)

[Louis, et al...Hypergraph Markov Operators, Eigenvalues and Approximation Algorithms STOC 2014](#)

[Chan, et al...Generalizing the Hypergraph Laplacian via a Diffusion Process with Mediators TCS 2020](#)

[Hayashi, et al...Hypergraph Random Walks, Laplacians, and Clustering CIKM 2020](#)

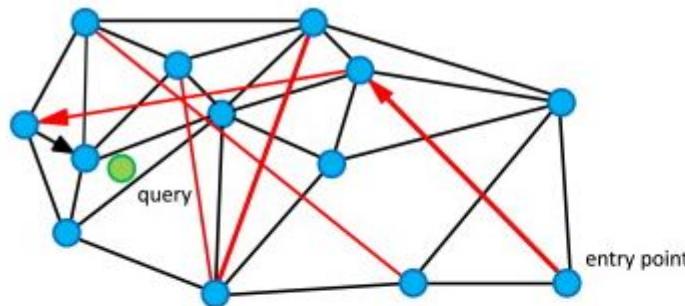
# Hypergraph NN



# Retrieval Method

2

# NSW



[Download](#) : [Download full-size image](#)

Fig. 1. Graph representation of the structure. Circles (vertices) are the data in metric space, black edges are the approximation of the Delaunay graph, and red edges are long range links for logarithmic scaling. Arrows show a sample path of the greedy algorithm from the entry point to the query (shown green). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

# Hierarchical NSW

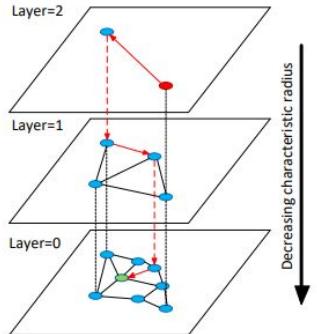


Fig. 1. Illustration of the Hierarchical NSW idea. The search starts from an element from the top layer (shown red). Red arrows show direction of the greedy algorithm from the entry point to the query (shown green).

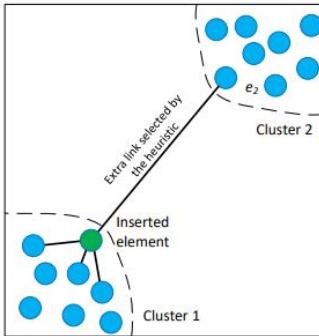
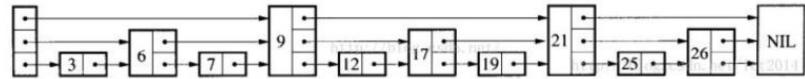
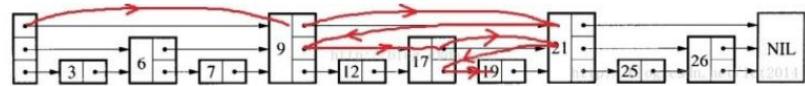


Fig. 2. Illustration of the heuristic used to select the graph neighbors for two isolated clusters. A new element is inserted on the boundary of Cluster 1. All of the closest neighbors of the element belong to the Cluster 1, thus missing the edges of Delaunay graph between the clusters. The heuristic, however, selects element  $e_2$  from Cluster 2, thus, maintaining the global connectivity in case the inserted element is the closest to  $e_2$  compared to any other element from Cluster 1.



查找

查找示意图如下:



比如我们要查找key为19的结点，那么我们不需要逐个遍历，而是按照如下步骤：

- 从header出发，从高到低的level进行查找，先索引到9这个结点，发现 $9 < 19$ ,继续查找(然后在level=2这层)，查找到21这个节点，由于 $21 > 19$ ,所以结点不往前走，而是level由2降低到1
- 然后索引到17这个节点，由于 $17 < 19$ ,所以继续往后，索引到21这个结点，发现 $21 > 19$ ,所以level由1降低到0
- 在结点17上，level==0索引到19,查找完毕。
- 如果在level=0这层没有查找到，那么说明不存在key为19的节点，查找失败

# Diffusion on Region Manifolds

- Contributions

-

# Aggregation

3

# ASMK (aggregated selective match kernels)

$$K(\mathcal{X}, \mathcal{Y}) = \gamma(\mathcal{X}) \gamma(\mathcal{Y}) \sum_{c \in \mathcal{C}} w_c M(\mathcal{X}_c, \mathcal{Y}_c)$$

$$\text{BoW} \quad M(\mathcal{X}_c, \mathcal{Y}_c) = \#\mathcal{X}_c \times \#\mathcal{Y}_c = \sum_{x \in \mathcal{X}_c} \sum_{y \in \mathcal{Y}_c} 1,$$

$$\text{BE} \quad M(\mathcal{X}_c, \mathcal{Y}_c) = \sum_{x \in \mathcal{X}_c} \sum_{y \in \mathcal{Y}_c} w(h(b_x, b_y))$$

$$\begin{aligned} \text{VLAD}^M(\mathcal{X}_c, \mathcal{Y}_c) &= V(\mathcal{X}_c)^\top V(\mathcal{Y}_c) \\ &= \sum_{x \in \mathcal{X}_c} \sum_{y \in \mathcal{Y}_c} r(x)^\top r(y). \end{aligned}$$

Non-aggregated kernel

$$M_N(\mathcal{X}_c, \mathcal{Y}_c) = \sum_{x \in \mathcal{X}_c} \sum_{y \in \mathcal{Y}_c} \sigma(\phi(x)^\top \phi(y)).$$

$$\text{SMK}(\mathcal{X}_c, \mathcal{Y}_c) = \sum_{x \in \mathcal{X}_c} \sum_{y \in \mathcal{Y}_c} \sigma_\alpha(\hat{r}(x)^\top \hat{r}(y)),$$

$$\sigma_\alpha(u) = \begin{cases} \text{sign}(u)|u|^\alpha & \text{if } u > \tau \\ 0 & \text{otherwise,} \end{cases} \quad \hat{r}(x) = \frac{x - q(x)}{\|x - q(x)\|}.$$

Aggregated kernel

$$M_A(\mathcal{X}_c, \mathcal{Y}_c) = \sigma \left\{ \psi \left( \sum_{x \in \mathcal{X}_c} \phi(x) \right)^\top \psi \left( \sum_{y \in \mathcal{Y}_c} \phi(y) \right) \right\} \quad (15)$$

$$= \sigma(\Phi(\mathcal{X}_c)^\top \Phi(\mathcal{Y}_c)), \quad (16)$$

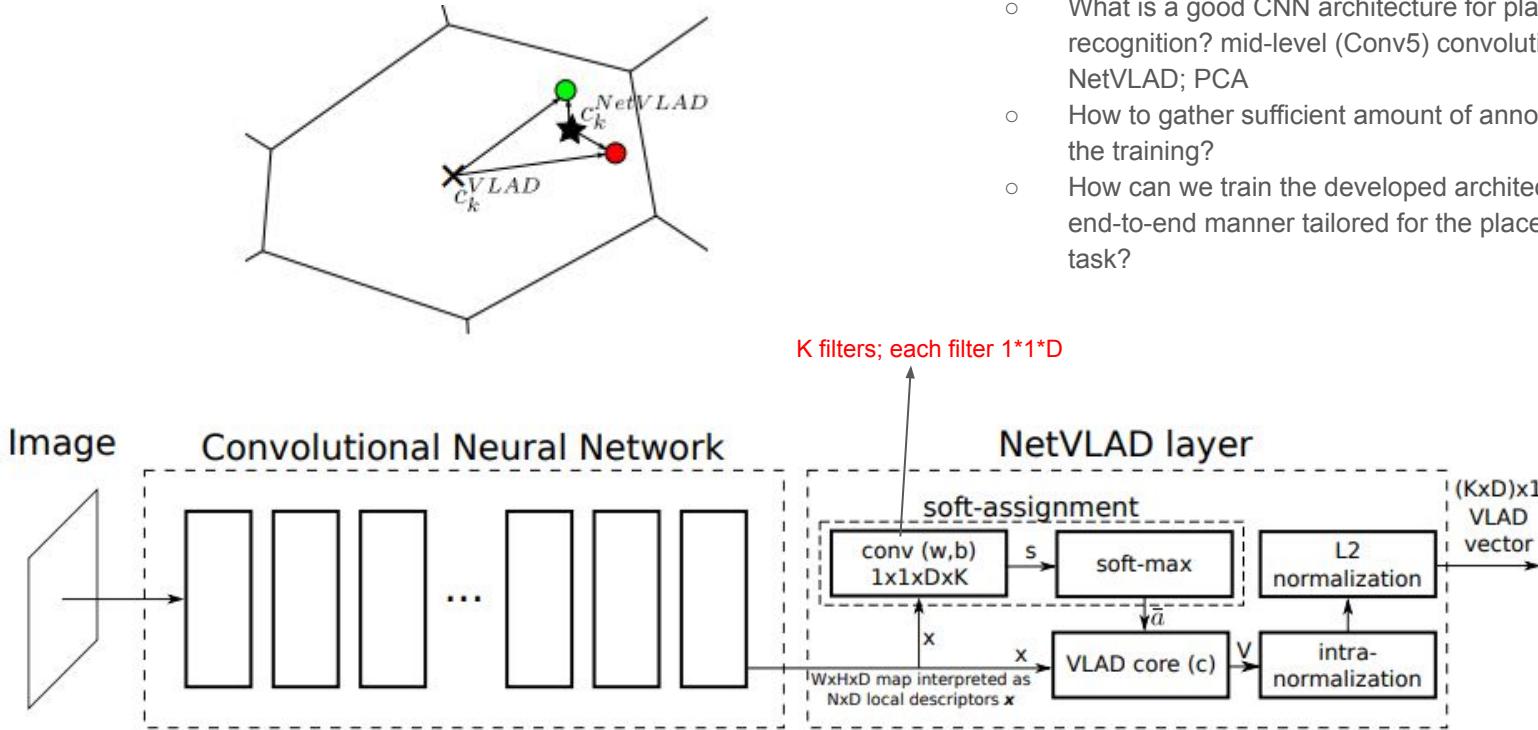
$$\text{ASMK}(\mathcal{X}_c, \mathcal{Y}_c) = \sigma_\alpha \left( \hat{V}(\mathcal{X}_c)^\top \hat{V}(\mathcal{Y}_c) \right)$$

$$\Phi(\mathcal{X}_c) = \hat{V}(\mathcal{X}_c) = V(\mathcal{X}_c) / \|V(\mathcal{X}_c)\|$$

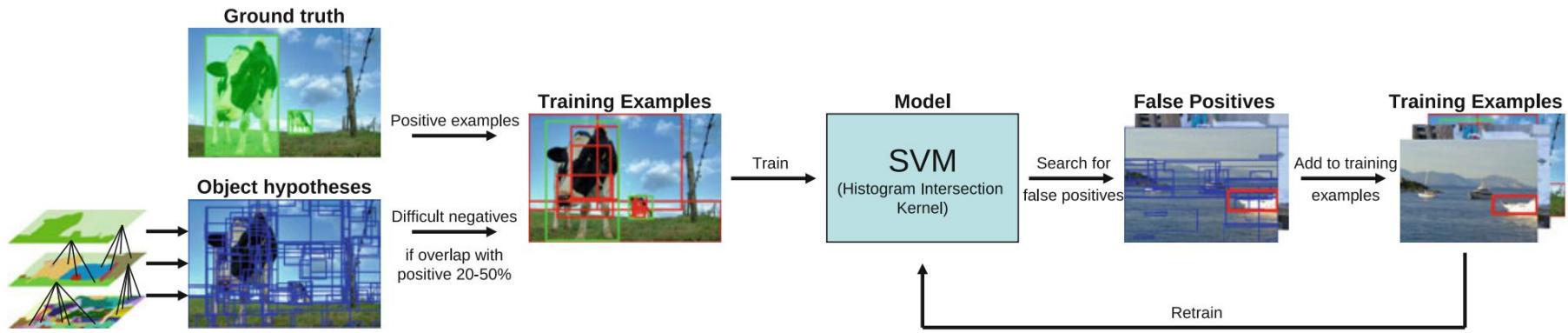
# Code NetVLAD

- Contributions

- What is a good CNN architecture for place recognition? mid-level (Conv5) convolutional features; NetVLAD; PCA
- How to gather sufficient amount of annotated data for the training?
- How can we train the developed architecture in an end-to-end manner tailored for the place recognition task?



# Selective Search



[OpenCV](#)

# Selective Search

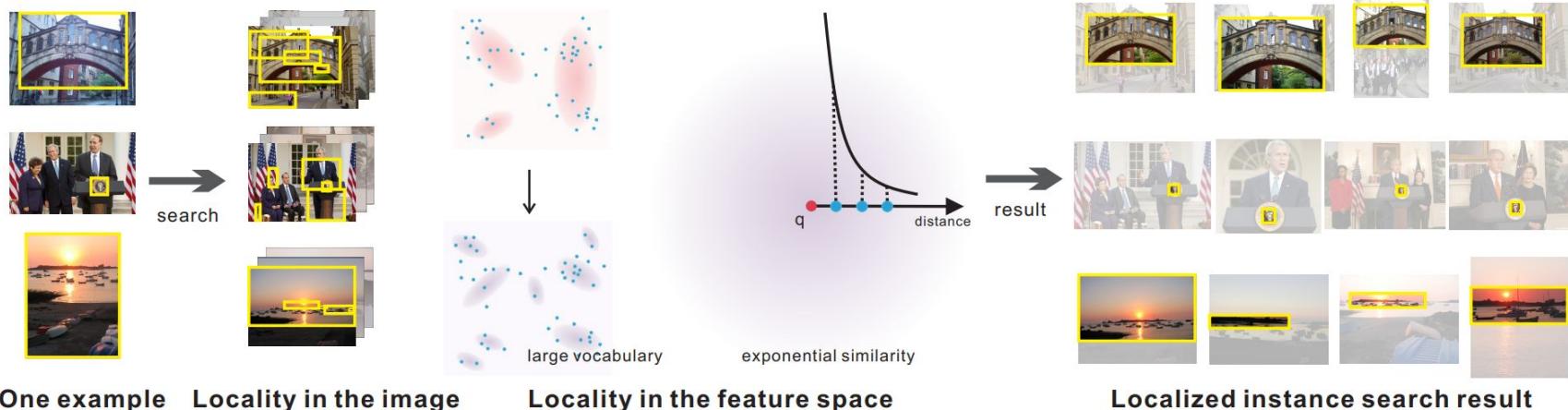


Figure 1. We propose locality in generic instance search from one example. As the first novelty, we consider many boxes as candidate targets to search locally in the picture by an efficient point-indexed representation. The same representation allows, as the second novelty, the application of very large vocabularies in Fisher vector and VLAD to search locally in the feature space. As the third novelty, we propose the exponential similarity to emphasize local matches in feature space. The method does not only improve the accuracy but also delivers a reliable localization.

# GeM (generalized-mean pooling)

- Contributions

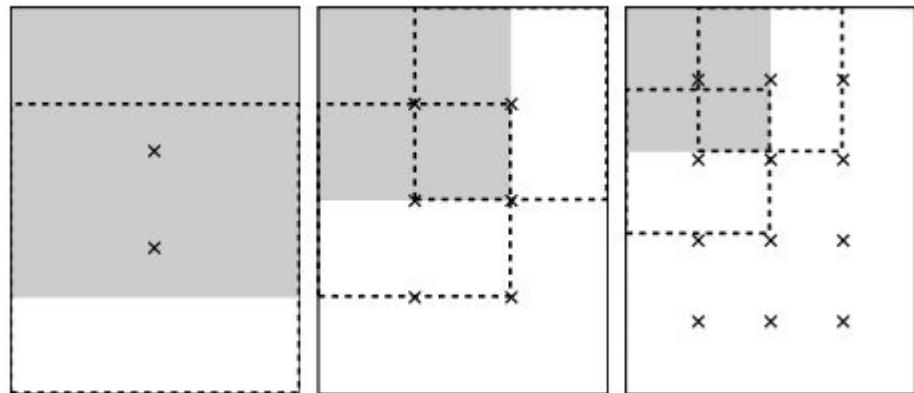
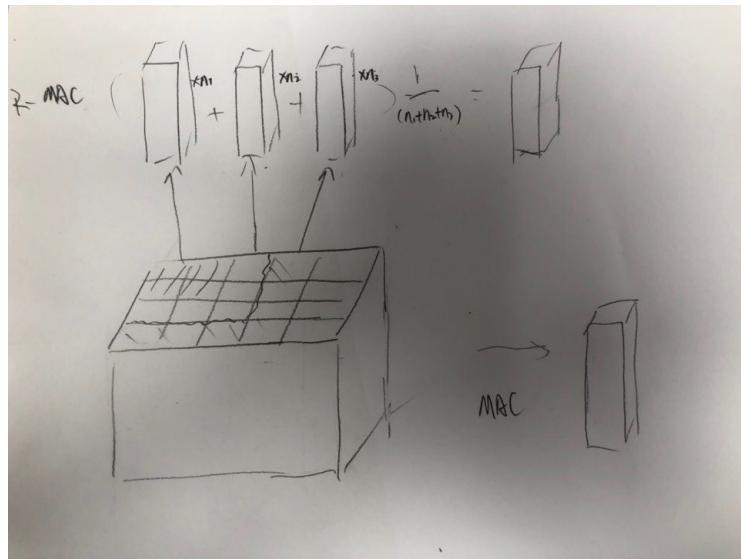
$$\mathbf{f}^{(g)} = [\mathbf{f}_1^{(g)} \dots \mathbf{f}_k^{(g)} \dots \mathbf{f}_K^{(g)}]^\top, \quad \mathbf{f}_k^{(g)} = \left( \frac{1}{|\mathcal{X}_k|} \sum_{x \in \mathcal{X}_k} x^{p_k} \right)^{\frac{1}{p_k}}.$$

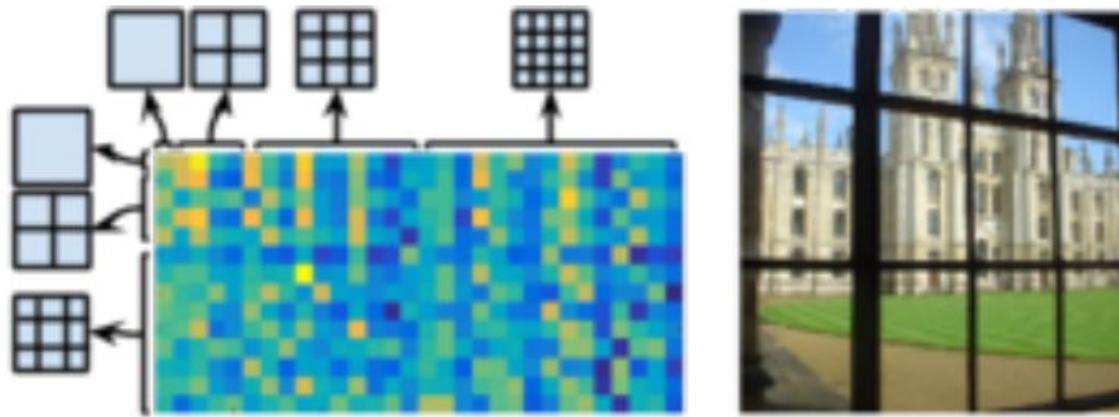
$$\frac{\partial \mathbf{f}_k}{\partial x_i} = \frac{1}{|\mathcal{X}_k|} \mathbf{f}_k^{1-p_k} x_i^{p_k-1},$$

$$\frac{\partial \mathbf{f}_k}{\partial p_k} = \frac{\mathbf{f}_k}{p_k^2} \left( \log \frac{|\mathcal{X}_k|}{\sum_{x \in \mathcal{X}_k} x^{p_k}} + p_k \frac{\sum_{x \in \mathcal{X}_k} x^{p_k} \log x}{\sum_{x \in \mathcal{X}_k} x^{p_k}} \right).$$

- Exploit **SfM information** and enforce, not only hard non-matching (negative), but also hard-matching (positive) examples from CNN training.
- Show that the whitening traditionally performed on short representations is, in some cases, unstable. So propose to learn the whitening through the same training data.
- Propose a **trainable pooling layer** that generalizes existing popular pooling schemes for CNNs.
- Propose a novel weighted query expansion that is more robust.
- A new state-of-the-art result.

# R-MAC

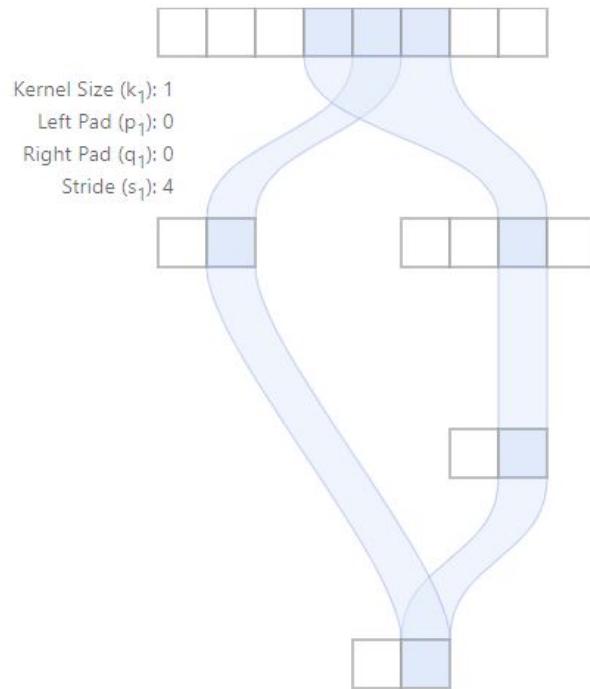




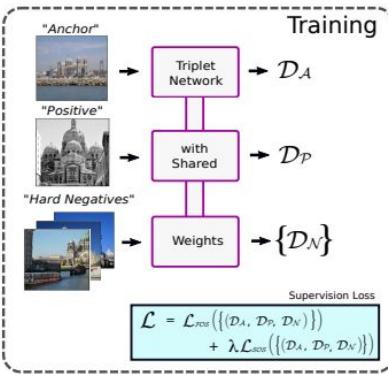
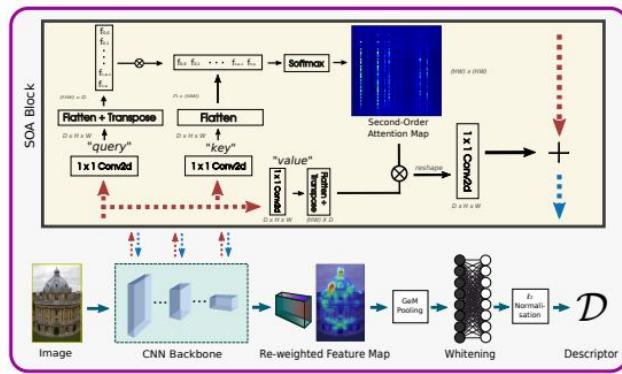
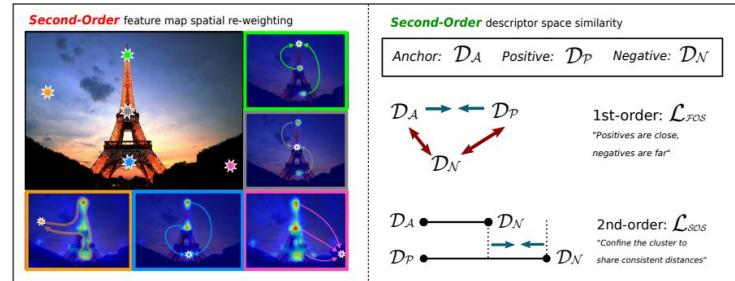
**Fig. 5 Schematic of distance matrix between two images:**

The distance between two images is computed by pooling the distances between different patches of the query and the reference image according to equations (3-4).

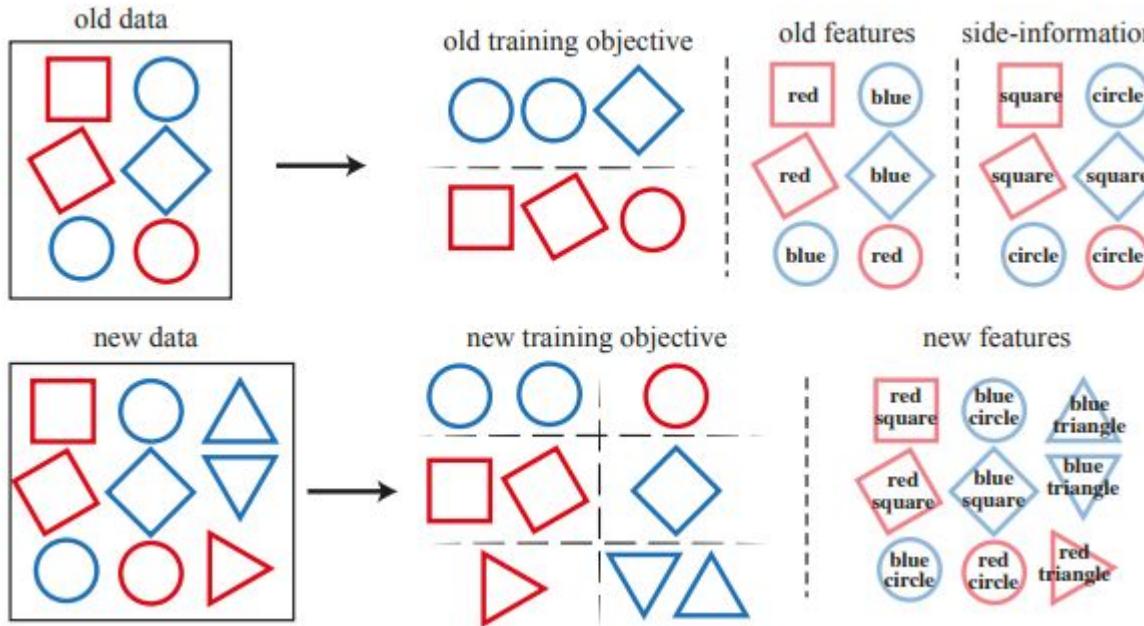
# Computing receptive fields



# Second-Order Loss



# Forward Compatible Training



# Whiten

# burstiness

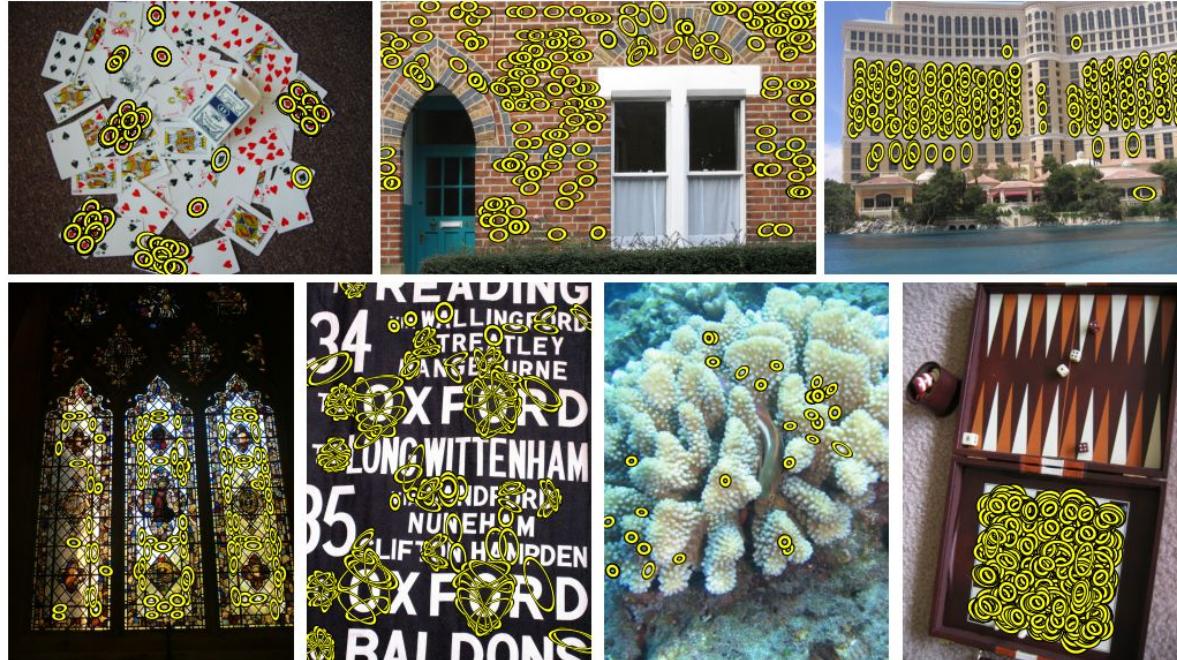


Figure 1. Illustration of burstiness. Features assigned to the most “bursty” visual word of each image are displayed.

# Co-occurrence

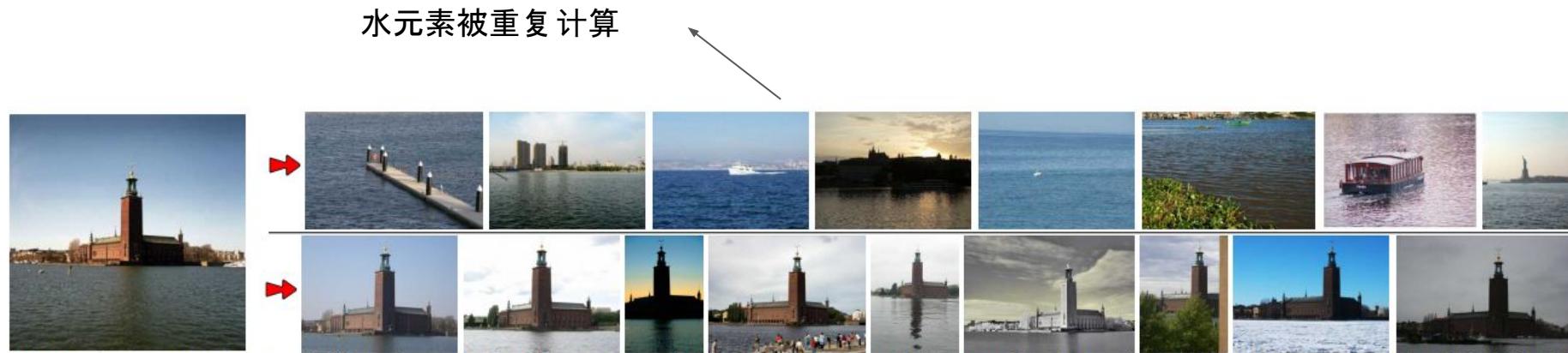


Figure 1. A query image (left) and top ranked images retrieved by two methods from a database of 5 million images: standard (*tf-idf* with spatial verification) retrieval (top) and the same method after automatic detection and removal of co-occurring features (bottom).

# Benefit of PCA and whitening: Co-missing and co-occurrence

$$\hat{X} = \frac{\text{diag}(\lambda_1^{-\frac{1}{2}}, \dots, \lambda_{D'}^{-\frac{1}{2}}) \mathbf{P}^\top X}{\left\| \text{diag}(\lambda_1^{-\frac{1}{2}}, \dots, \lambda_{D'}^{-\frac{1}{2}}) \mathbf{P}^\top X \right\|},$$

# Classical Features

# Hessian-affine detector

[Perdoch, et al... Efficient repersentation of local geometry for image scale object retrieval, CVPR 2009](#)

# rootSIFT descriptor

[Arandjelovic, et al... Three things everyone should know to improve object retrieval, CVPR 2012](#)

# CNN Features

5

# Faster R-CNN features for Instance Search

- Global Features + Local Features

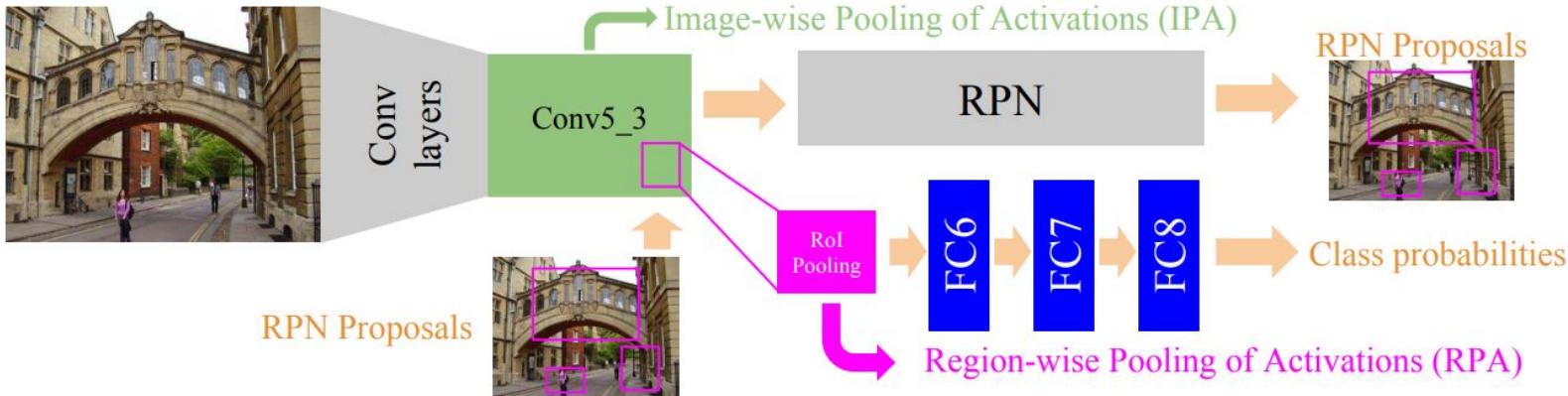


Figure 2. Image- and region-wise descriptor pooling from the Faster R-CNN architecture.

# Bags of Convolutional Features

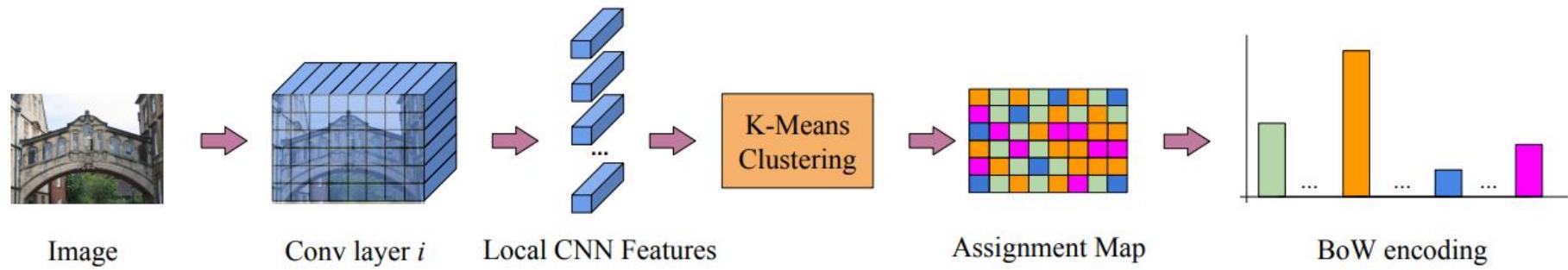
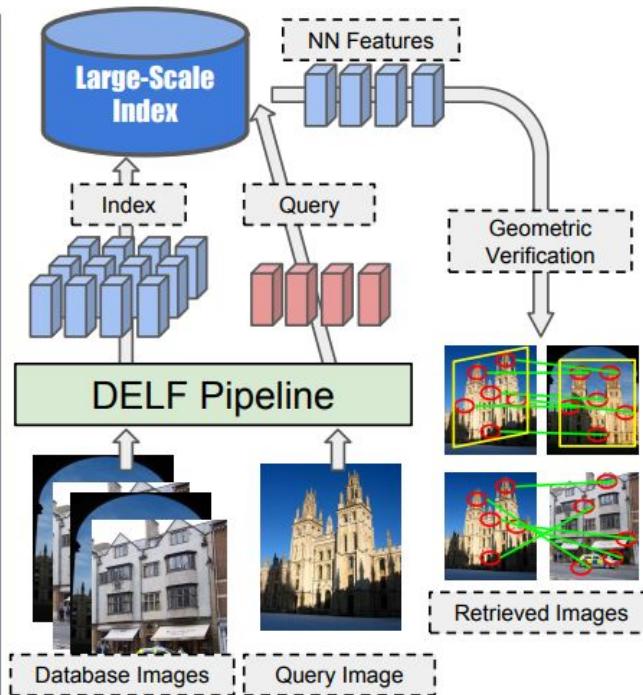
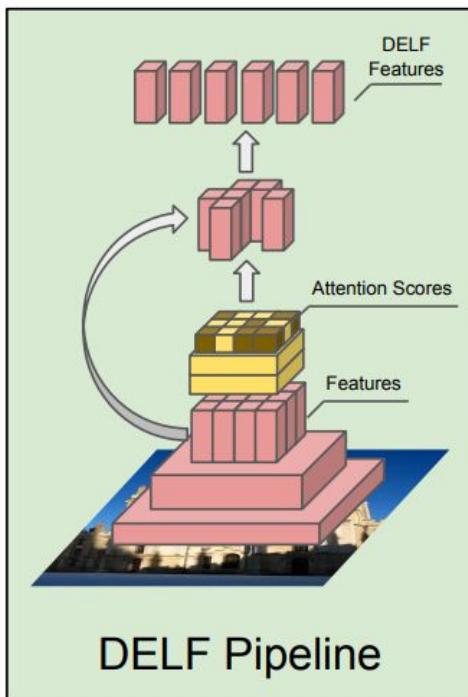


Figure 2: The Bag of Local Convolutional Features pipeline (BLCF).

# DELF (deep local features)



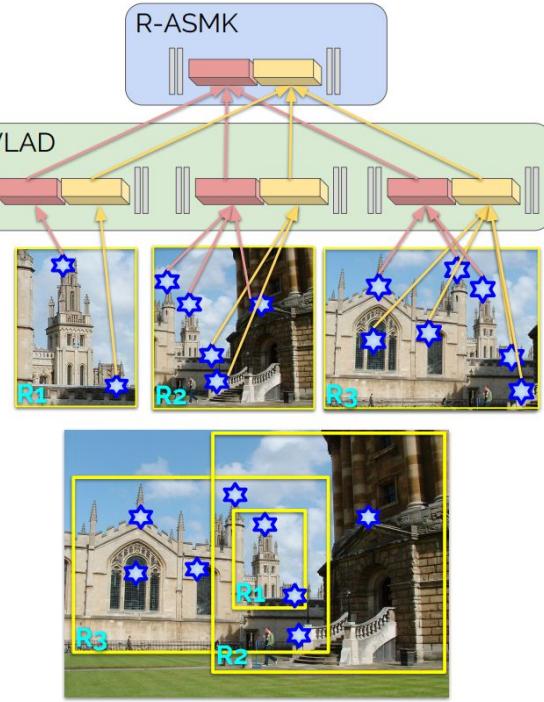
- Contributions

- A new large-scale dataset, Google-Landmarks, which contains more than 1 M landmark images from almost 13K unique landmarks.
- Propose DELF, a CNN-based local feature with attention, which is trained with weak supervision using image-level class labels only.

# Code

## Detect-to-Retrieve

Regional Aggregation



Local Feature Extraction  
+ Region Detection

- Contributions

- Propose Google Landmark Boxes datasets (add landmark boxes on Google Landmark dataset).
- Leverage the trained detector and produce more efficient regional search systems.
- Propose regional aggregated match kernels to leverage selected image regions and produce a discriminative image representation.

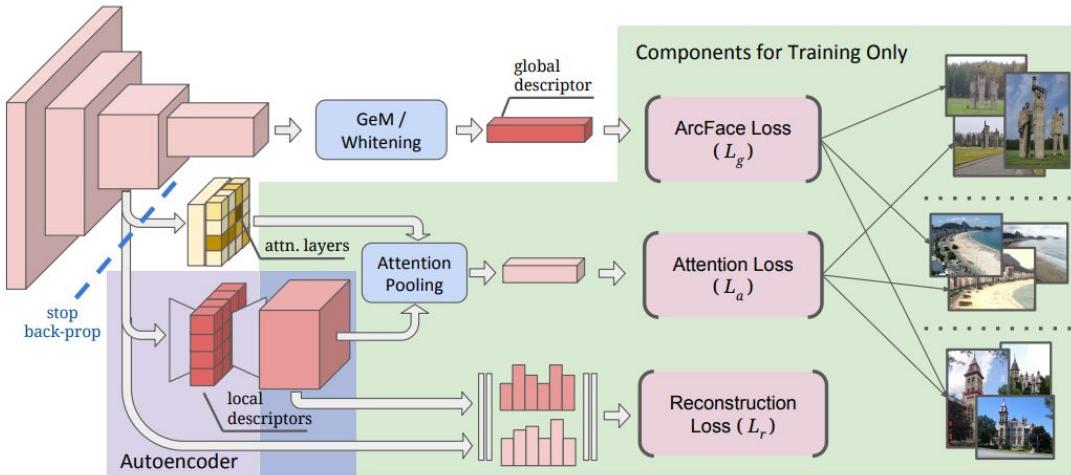
$$\begin{aligned} \text{sim}^{(\text{R-VLAD})}(X, Y^{(n)}) &= \gamma(\mathcal{X}) \sum_c \sum_r \frac{\gamma(\mathcal{Y}_c^{(n,r)})}{R_n} V(\mathcal{X}_c)^T V(\mathcal{Y}_c^{(n,r)}) \\ &= \sum_c \gamma(\mathcal{X}) V(\mathcal{X}_c)^T \sum_r \frac{\gamma(\mathcal{Y}_c^{(n,r)})}{R_n} V(\mathcal{Y}_c^{(n,r)}) \\ &= \sum_c V_R(\mathcal{X}_c)^T V_R(\{\mathcal{Y}_c^{(n,r)}\}_r) \end{aligned}$$

$$\text{sim}^{(\text{R-ASMK})}(X, Y^{(n)}) = \sum_c \sigma \left( \hat{V}_R(\mathcal{X}_c)^T \hat{V}_R(\{\mathcal{Y}_c^{(n,r)}\}_r) \right)$$

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

出现次数较少的visual patterns  
应该更重要

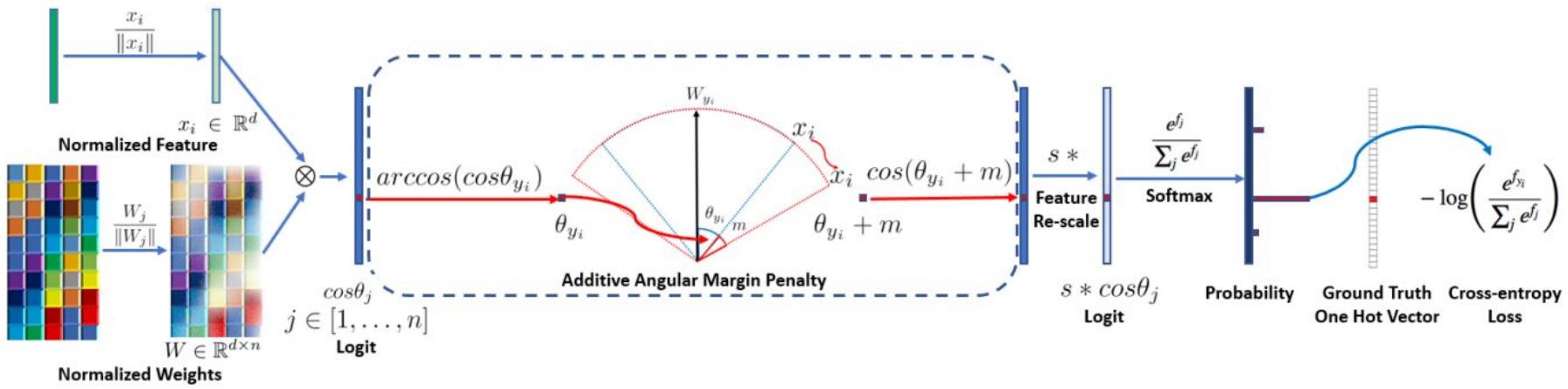
# Code DELG



- Contributions

- **DELG** (DEep Local and Global features), a unified model to represent both local and global features, using a convolutional neural network (CNN).
- Adopt a convolutional **autoencoder module** that can successfully learn low-dimensional local descriptors. (avoid the need of post-processing learning steps, such as PCA.)
- Design a procedure that enables end-to-end training of the proposed model using only image-level supervision. (**control the gradient flow**)

# ArcFace



# Unsupervised

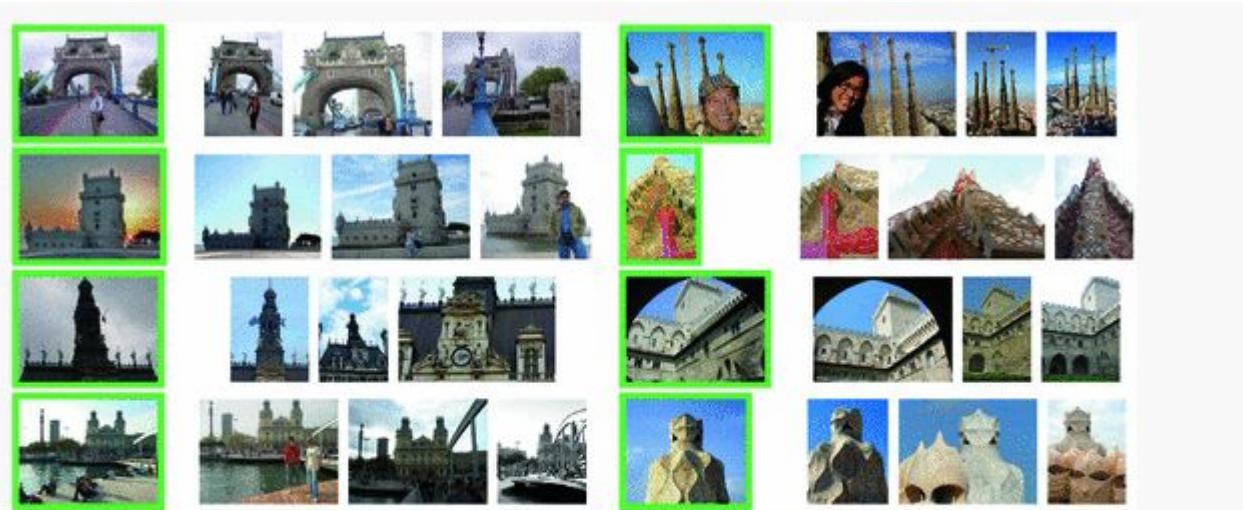
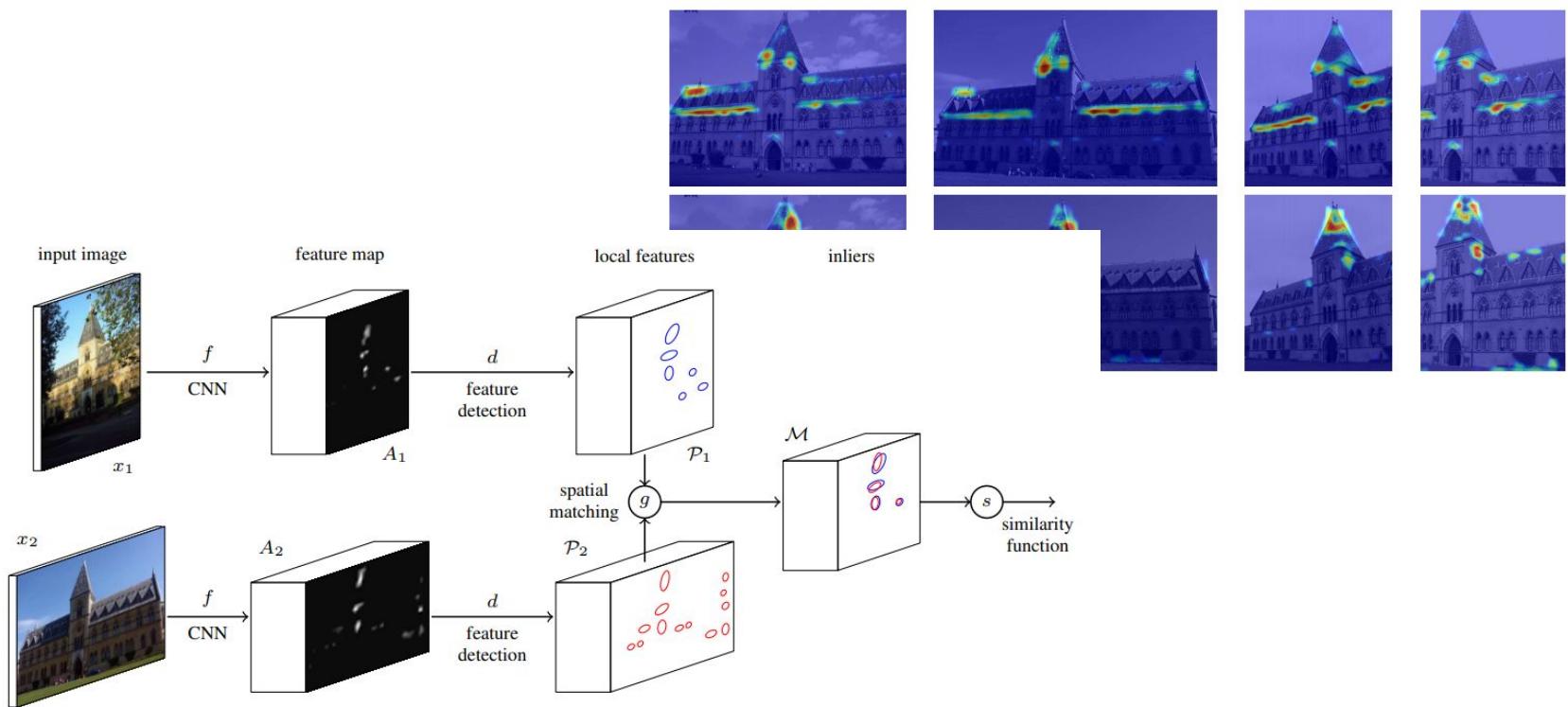


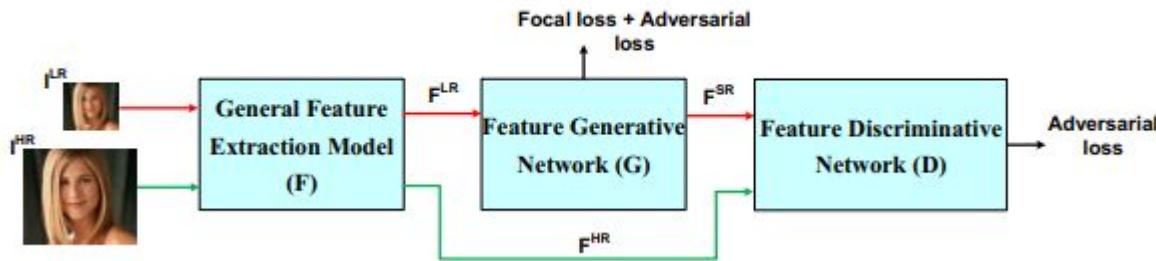
Fig. 2.

Examples of training query images (green border) and matching images selected as positive examples by methods (from left to right)  $m_1(q)$ ,  $m_2(q)$ , and  $m_3(q)$ . (Color figure online)

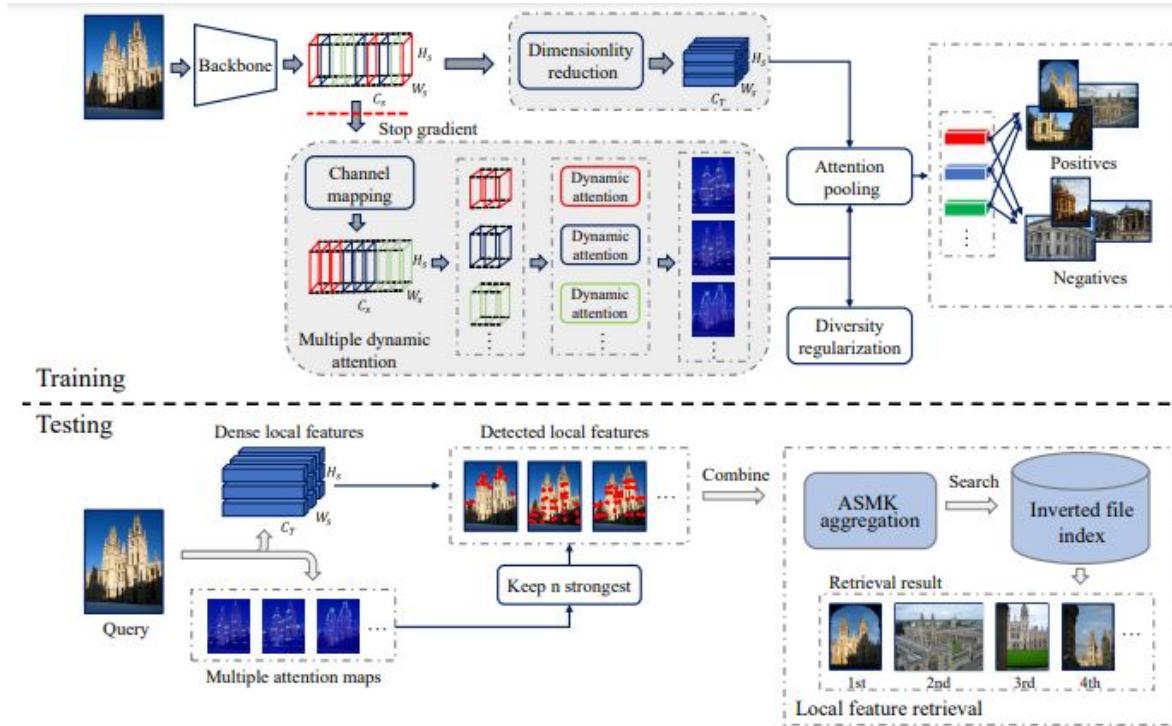
# Spatial matching per channel



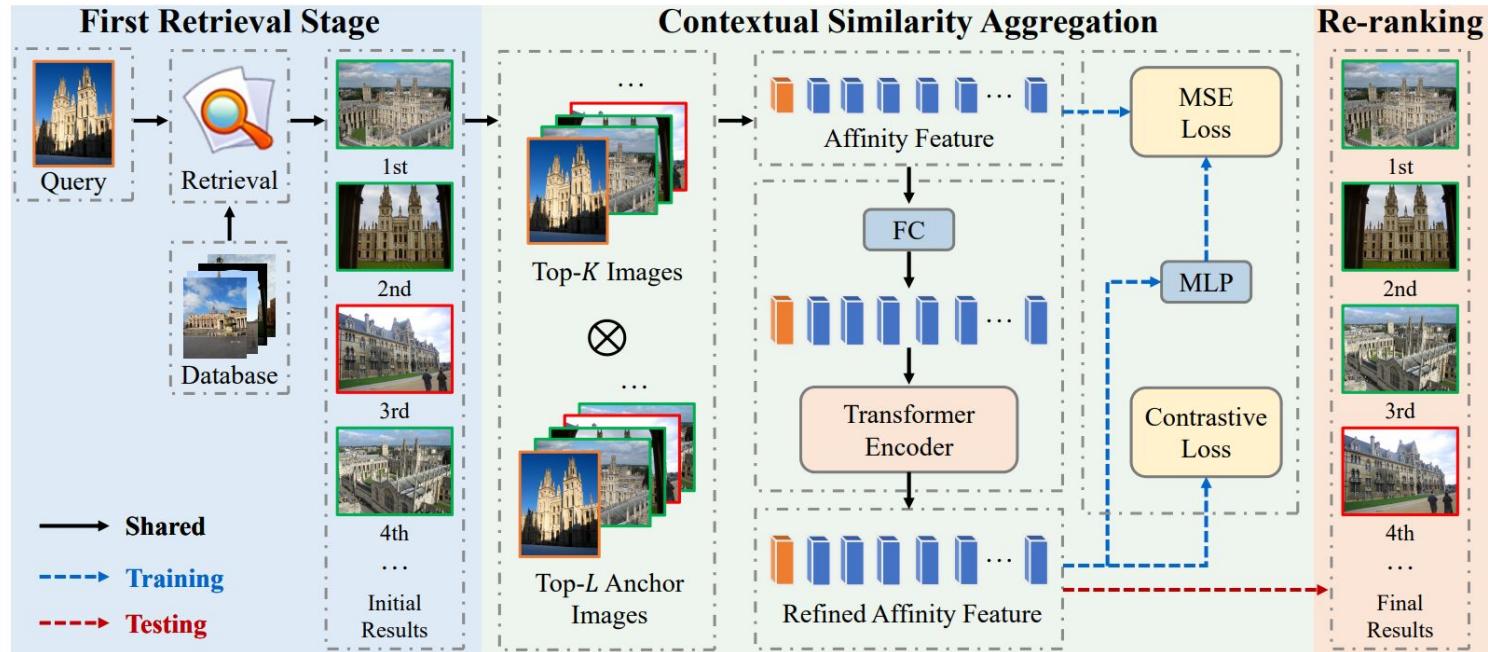
# Feature Super-Resolution



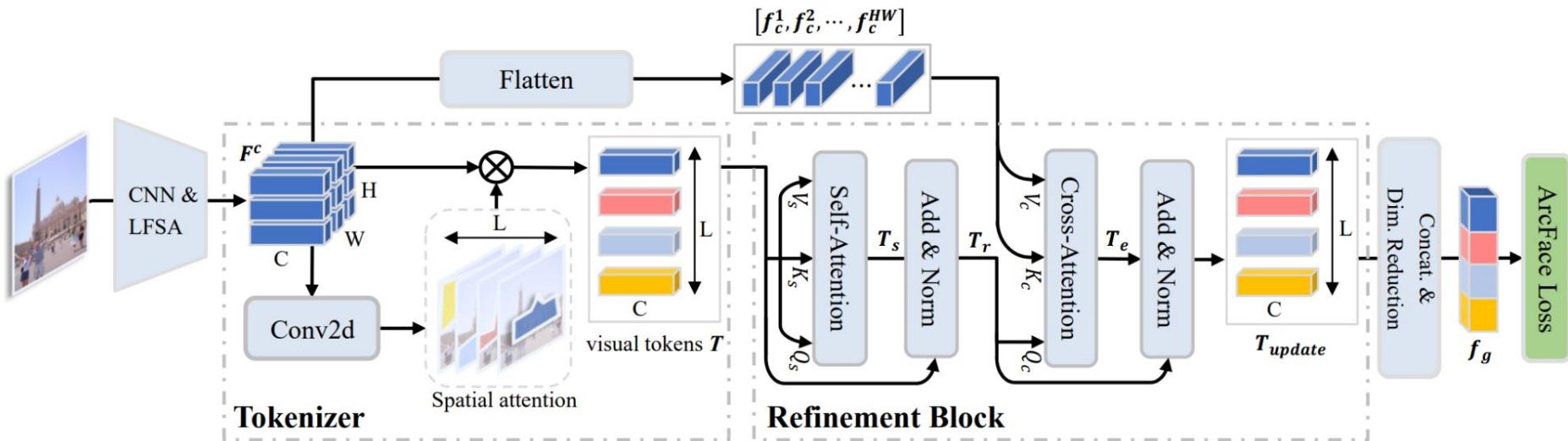
# Attention local features



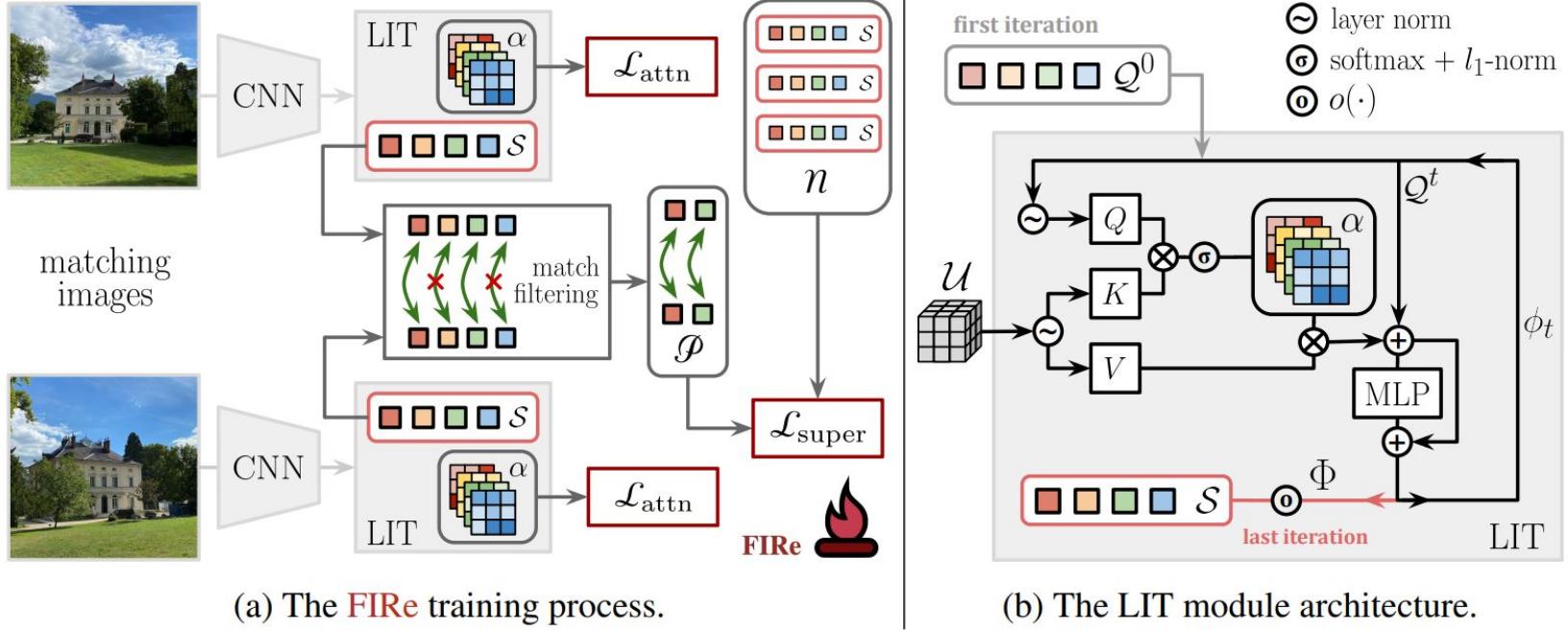
# Transformer + query expansion



# Transformer + query expansion



# Contrastive learning + transformer



# Sub-image search

2

# Object-based image retrieval

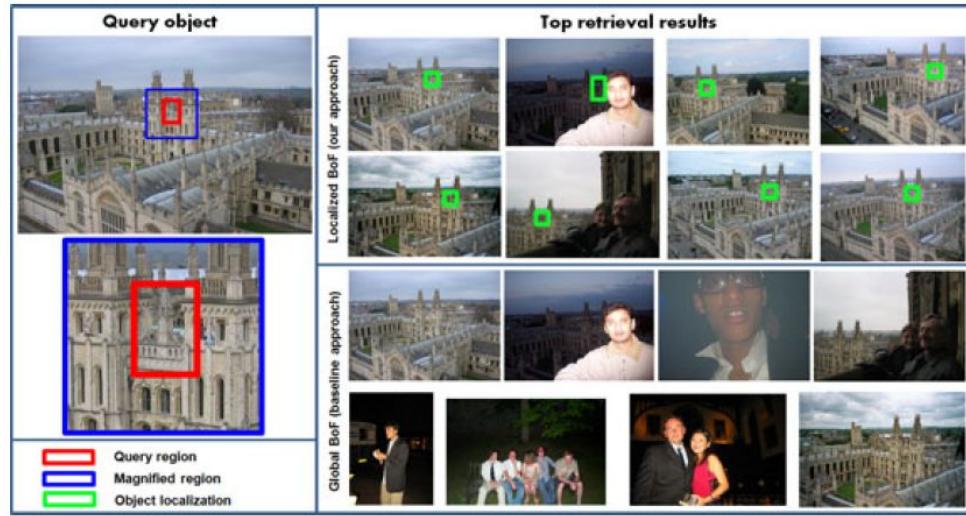
- The combination of content based image retrieval + object localization
  - Give a object patch, how to find the same object in the database.
  - High accuracy + large set of candidate images.



Figure 1. *Object-based image retrieval.* Left: The users marks an object in an image by a bounding region. Right: The system returns images from a database that show the query object.

# A local Bag-of-Features Model

Represent each database images as a family of histograms that depend functionally on a bounding rectangle.



**Fig. 1.** An example of small object retrieval. Left: a query image with a region of interest. Right: the top 8 retrieved images using our approach and the baseline Global BoF approach. Red rectangle represents the query object of interest, and Green rectangles represent the returned object bounding boxes using our approach. In contrast, the baseline method cannot return bounding boxes and need additional totally different criterion (*e.g.* RANSAC) to localize objects.

# Unsupervised Object discovery

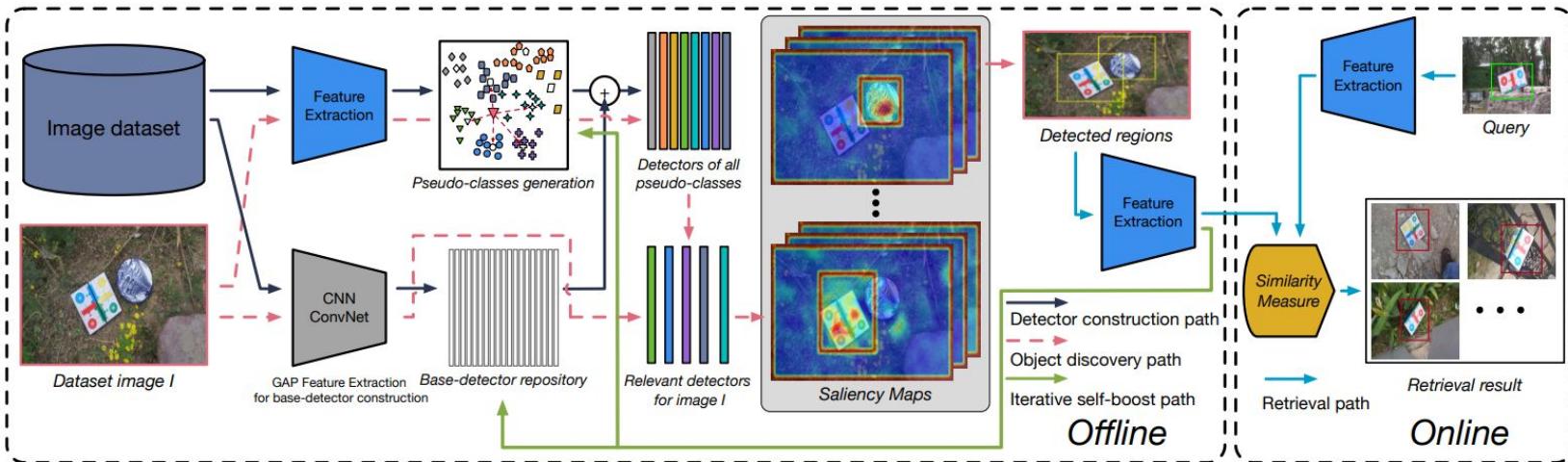


Fig. 1. Overview of the proposed dataset-driven unsupervised object discovery framework (DUODIS in short) for region-based instance image retrieval. Our framework consists of two (i.e., offline and online) parts. In the offline part, there are three kinds of paths. The detector construction path (in black) exploits the whole image dataset and constructs the base-detector repository and the object detectors for all pseudo-classes. The object discovery path (in red and dashed line) assigns relevant object detectors to each image and applies them to conduct object discovery. The iterative self-boosting path (in green) utilizes the detected regions to further enhance the whole framework. In the online part, once a query object is submitted, its features will be extracted to measure its similarity to each of the regions detected from the image dataset. The most similar images are then retrieved, with the matched object regions delineated. In instance image retrieval, it is usually assumed that a user will indicate the object to search for in a query image by a bounding box. This setting is followed in this paper.

# Grid-based



**Fig. 1.** An example of visual object search. Left: a query object, such as a logo, selected by the user. Right: resulting images from visual object search, where object locations are identified and marked by blue bounding boxes.

# Branch and Bound

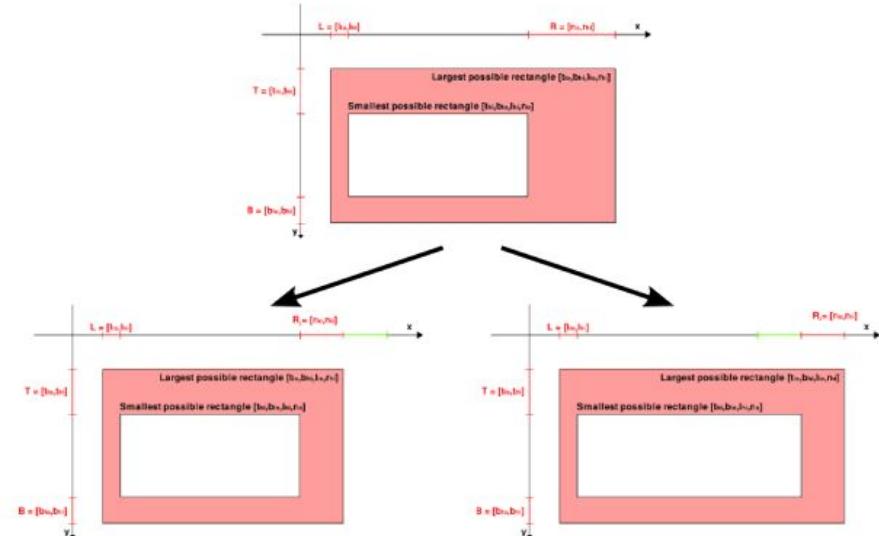


Fig. 2. Splitting rectangle sets is done by dividing one of the intervals in two. In this case,  $[T, B, L, R] \rightarrow [T, B, L, R_1] \cup [T, B, L, R_2]$ , where  $R_1 := [r_{lo}, \lfloor \frac{r_{lo}+r_{hi}}{2} \rfloor]$  and  $R_2 := [\lfloor \frac{r_{lo}+r_{hi}}{2} \rfloor + 1, r_{hi}]$ .



Fig. 3. Example images of cat (top) and dog (bottom) categories of PASCAL VOC 2006 dataset. Objects occur in different sizes and poses, and multiple object instances are possible within one image. Objects are also frequently occluded or truncated.

[Lampert, et al., Efficient Subwindow Search: A Branch and Bound Framework for Object Localization, TPAMI 2009](#)

# Strain estimation



Fig 2: Specimen fixed in UTM

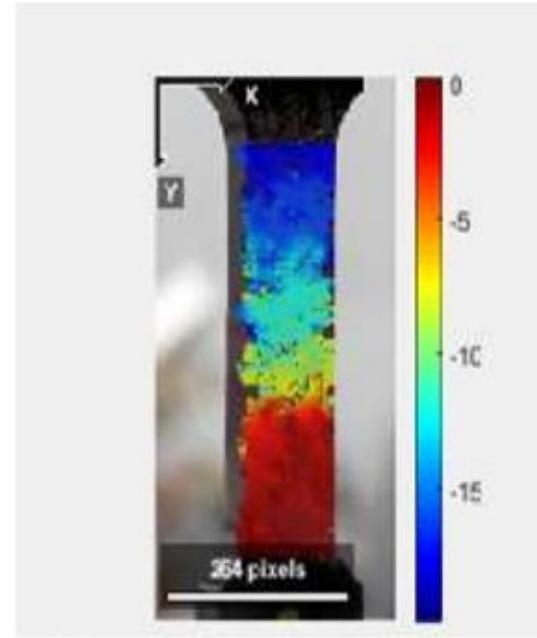


Fig. 5. Strain maps curves obtained on aluminium specimen

digital pathology require sub-image retrieval rather than the whole image retrieval for the system to be of **clinical use**

**Abstract:**

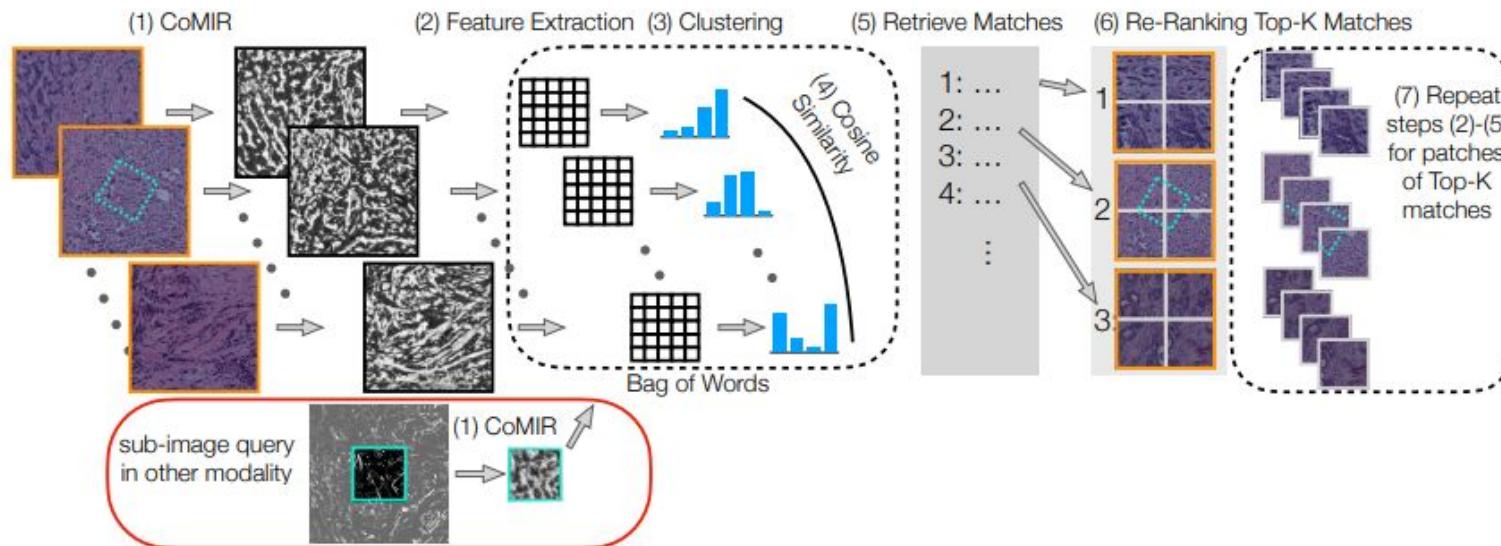
Content-based image retrieval systems for digital pathology require sub-image retrieval rather than the whole image retrieval for the system to be of clinical use.

Digital pathology images are huge in size and thus the pathologist is interested in retrieving specific structures from the whole images in the database along with the previous diagnosis of the retrieved sub-image. We propose a content-based sub-image retrieval system (sCBIR) framework for high resolution digital pathology images. We utilize scale-invariant feature extraction and present an efficient and robust searching mechanism for indexing the images as well as for query execution of sub-image retrieval. We present a working sCBIR system and show results of testing our system on a set of queries for specific structures of interest for pathologists in clinical use. The outcomes of the sCBIR system are compared to manual search and there is an 80% match in the top five searches.

[Eva Breznik et al...Content based sub-image retrieval system for high resolution pathology images using salient interest points](#)

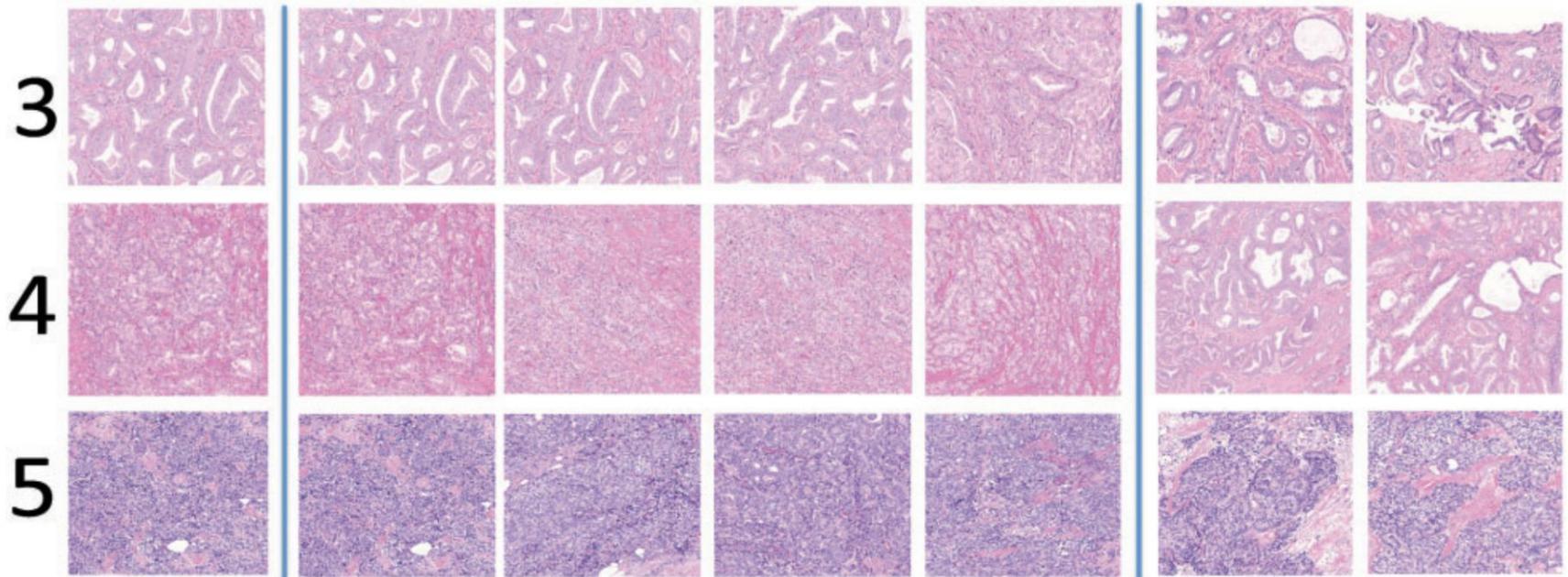
[2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society](#)

# improve diagnosis and discover patterns in pathologies



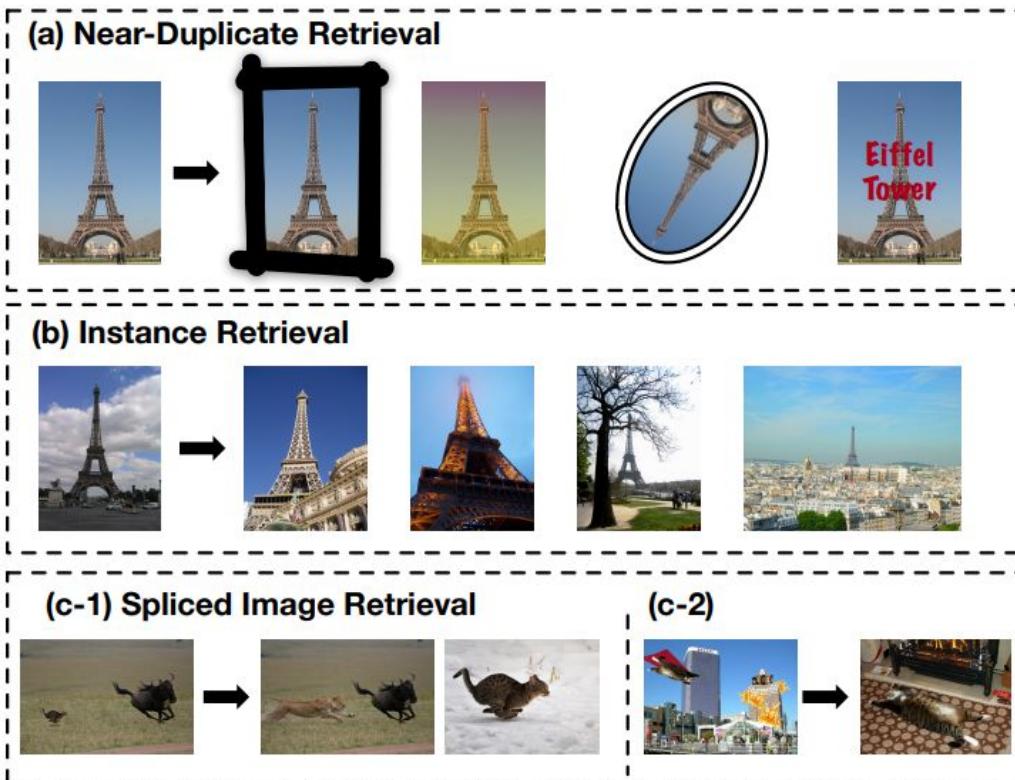
**Figure 1:** The proposed s-CBIR pipeline learns CoMIRs for the images in the repository and the query in form of a patch, extracts sparse features, which are binned into single descriptors for each image, building the vocabulary for a BoW. Matches are found using the cosine similarity. The Top-K matches are split into patches and a new BoW is computed for Re-Ranking.

# Bio

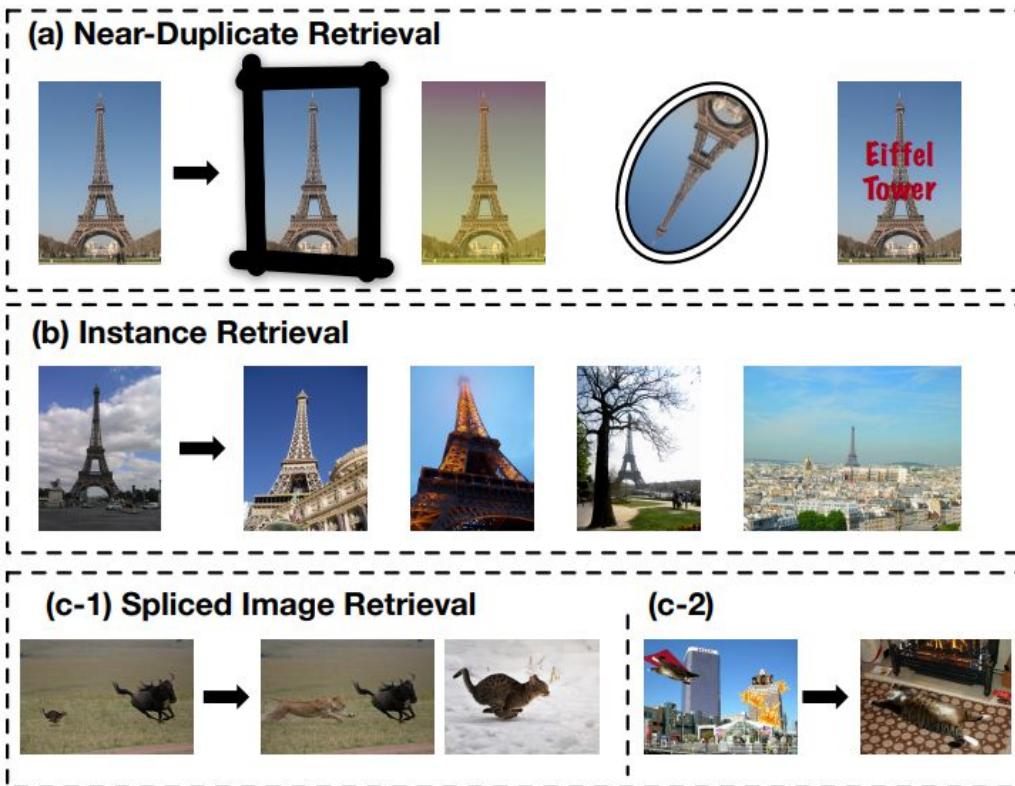


**Fig. 4.** The sub-image content-based rank retrieval results. The left panel is the query image; the middle panel is the retrieval results representing the most similar cases. The right panel is the retrieval results representing the most dissimilar cases. The number 3, 4 and 5 correspond to different Gleason scores

# Spliced Image Retrieval



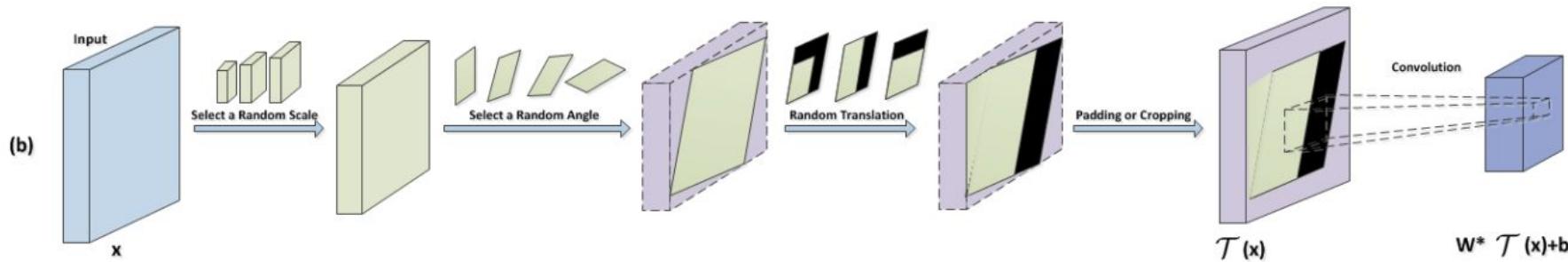
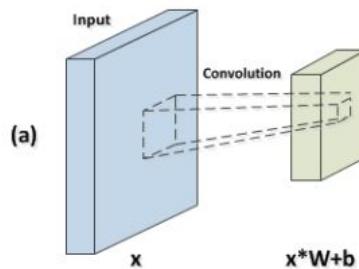
# Spliced Image Retrieval



# Transform-Invariant

1

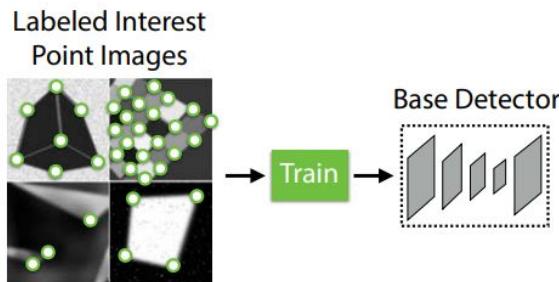
# Transformed Invariant Search



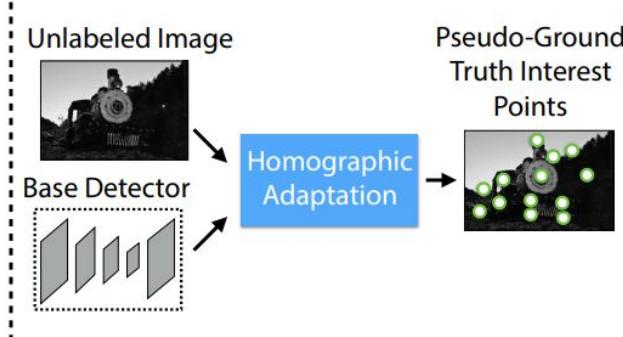
# Correspondence features & Interest Point

# SuperPoint

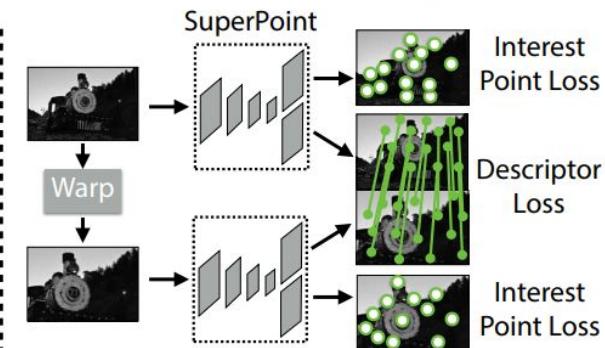
(a) Interest Point Pre-Training



(b) Interest Point Self-Labeling



(c) Joint Training



[see Section 4]

[see Section 5]

[see Section 3]

# Dense Pixel Matching

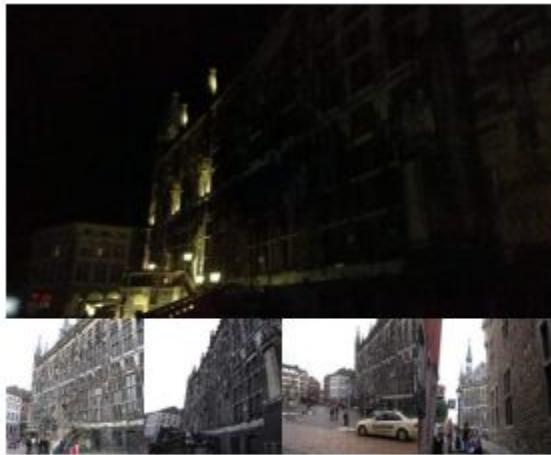


Figure 1: Qualitative results of the proposed method for the task of image retrieval. The first row is a query taken at night-time with a mobile phone camera and the last row is a list of top-4 retrieved database images obtained by our method. All 4 are correct matches.

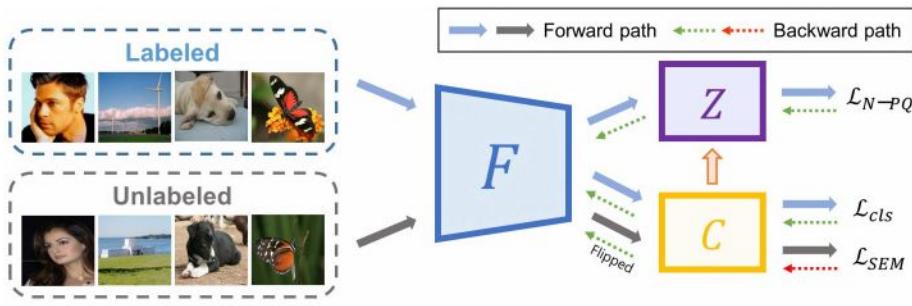


# Data Compression

1

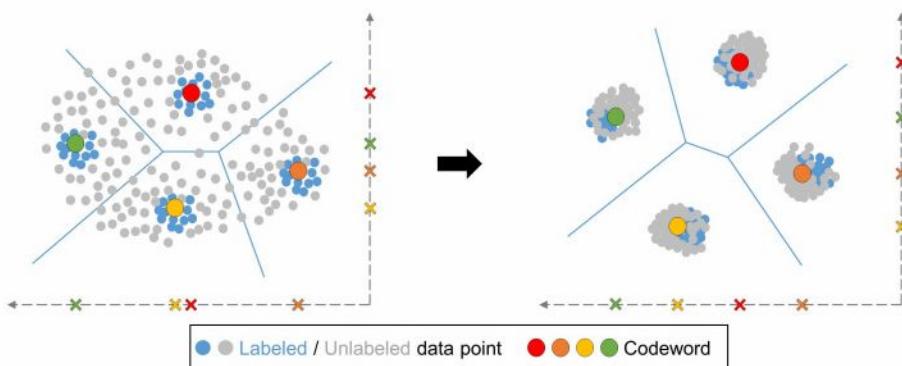
# Code

# GPQ (Generalized Product Quantization)



- Contributions

- The first deep semi-supervised PQ scheme for image retrieval.
- With the proposed **metric learning** strategy and **entropy regularization term**, the semantic similarity of labeled data is well preserved into the codewords, and the underlying structure of unlabeled data can be fully used to generalize the network.
- Yield the state-of-art retrieval results in semi-supervised image retrieval protocols.



# Feature Visualization

[pytorch\\_cnn\\_visualization](#)

# Visualizing and Understanding Convolutional Networks

Code

# Feature visualization

# Code

## Computing receptive fields

# Correlation Filters

0

# With RL; RNN

# Image clustering

2

# simHash for identical image detection

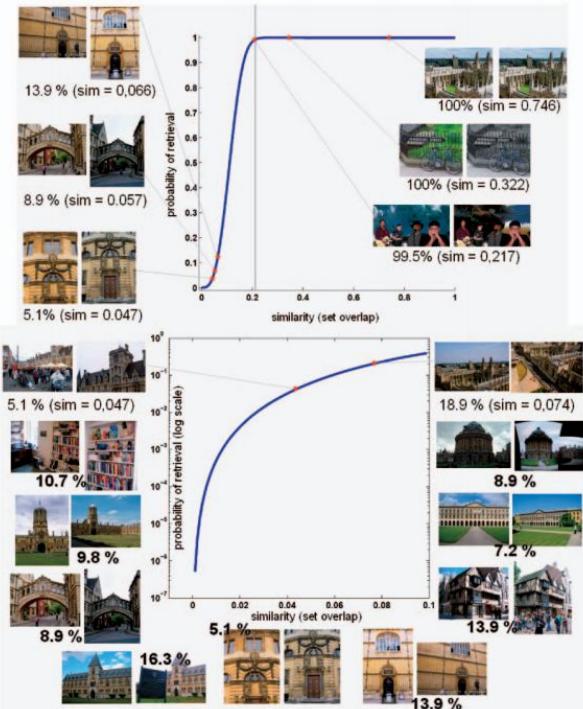
permutations	A	B	C	D	E	F	A B C	B C D	A E F
	3	6	2	5	4	1	2	2	1
	1	2	6	3	5	4	1	2	1
	3	2	1	6	4	5	1	1	3
	4	3	5	6	1	2	3	3	1

min-hashes

$$P(\min \pi(\mathcal{A}_1) = \min \pi(\mathcal{A}_2)) = \frac{|\mathcal{A}_1 \cap \mathcal{A}_2|}{|\mathcal{A}_1 \cup \mathcal{A}_2|} = \text{sim}(\mathcal{A}_1, \mathcal{A}_2).$$

$$P(\mathcal{A}_1 \approx \mathcal{A}_2) = \sum_{i=m}^k \binom{k}{i} p^{in} (1 - p^n)^{k-i},$$

# min-Hash Image Clustering (MHIC)



- 1) **Hashing.** Image descriptors are stored in a hash table. In our experiments,  $2^{51}$  different descriptor values are used. The probability of two images falling into the same bin (*exact* descriptor match) is proportional to their similarity – eqn. (2).
- 2) **Similarity estimation.** For all  $\binom{n}{2}$  pairs of the  $n$  images hashed in the same bin, i.e. for  $n$  collisions in the bin, similarity is calculated. The process is fast and consists of comparing two vectors and counting the number of identical elements. In our experiments, the number of vector elements is 512. The similarity is then thresholded.
- 3) **Spatial consistency.** For each image pair that passed the similarity test, spatial consistency is verified. Image pairs that pass the spatial consistency test become *cluster seeds*.
- 4) **Seed growing.** Seed images are used as visual queries and the query expansion technique is used to ‘crawl’ the images in the cluster.

Fig. 2. The min-Hash Image Clustering (MHIC) algorithm.





# Attack

# You See What I Want You to See

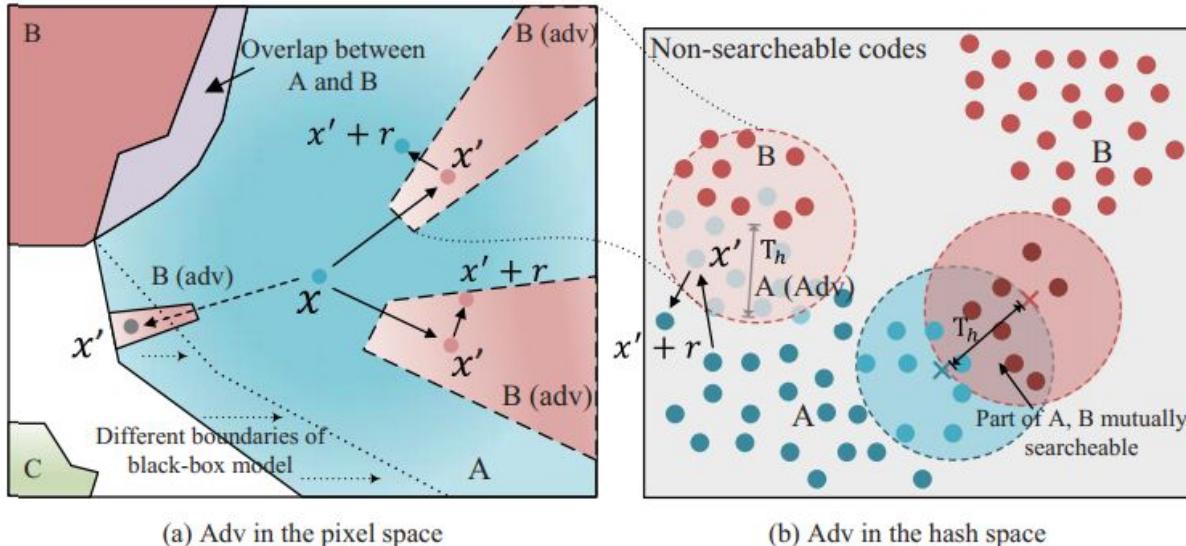
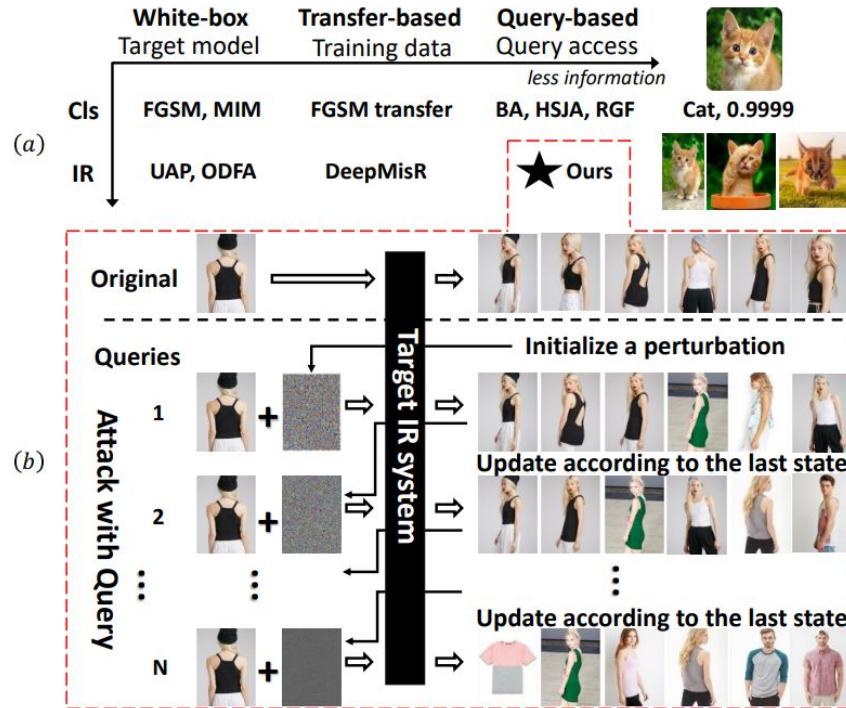


Figure 4: Illustration of adversarial examples (a) image pixel space (b) hash space.

# Query-efficient attacks



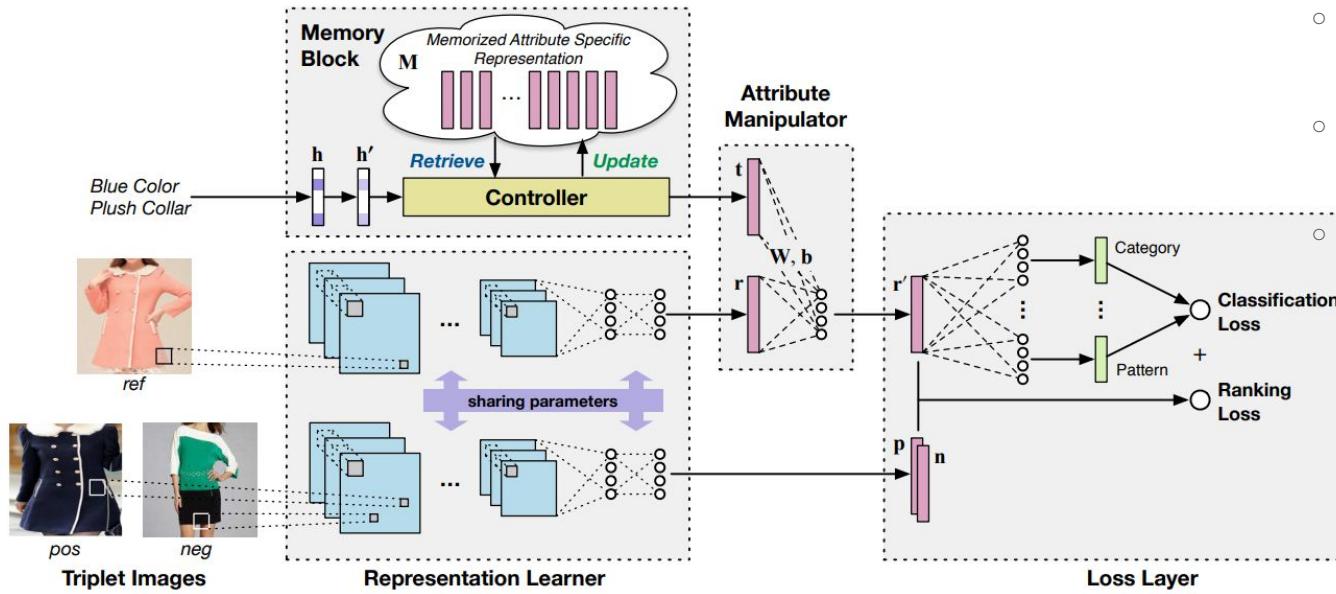
# Fashion Search

2

# Memory-augmented; Interactive Fashion Search

- Contributions

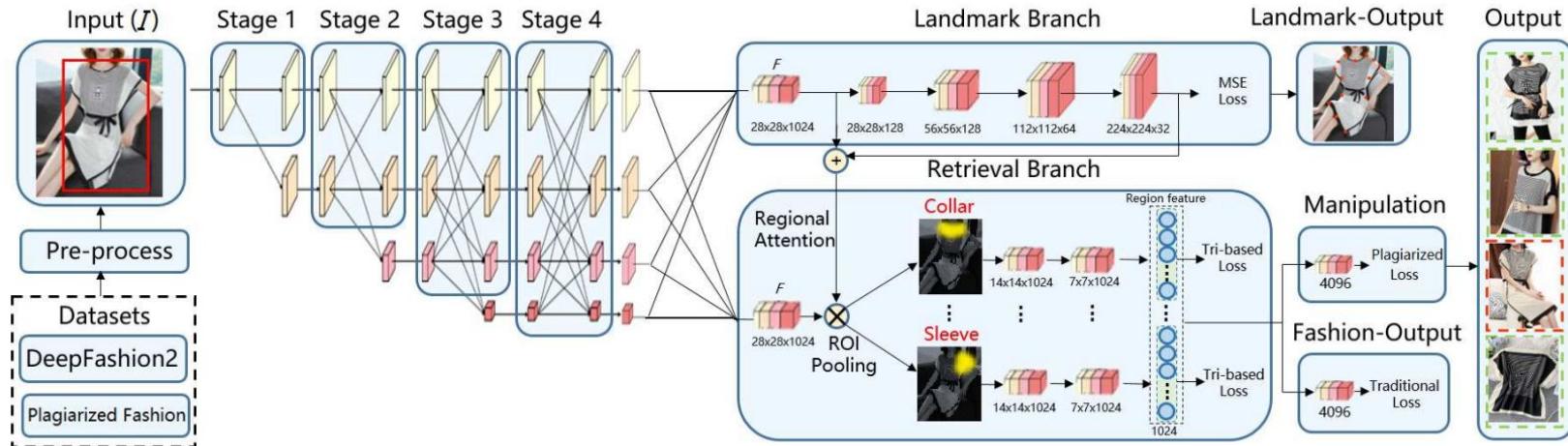
- Propose a novel memory-augmented attribute manipulation network for interactive fashion search.
- A memory block consisting of a memory and a neural controller in the proposed networks.
- Exploit a joint optimization method for attribute representation and retrieval feature learning.



# Plagiarized-Search-Net (PS-Net)

- Contributions

- Proposed PS-Net, which obtain the regional representation of images and compute the similarity region by region
- A novel problem and a dataset.
- PS-Net can also be used for traditional fashion retrieval and landmark estimation tasks.



# Content-style + text feedback



(a) Example of a **content** modification



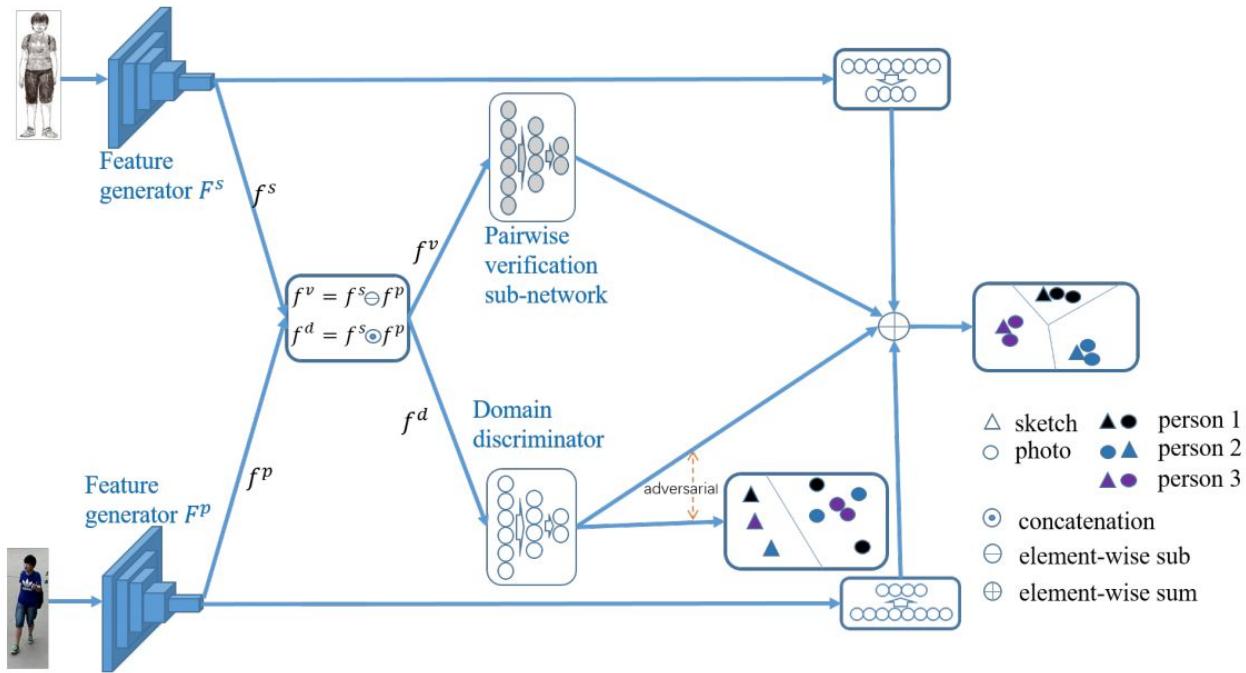
(b) Example of a **style** modification

# Sketch; Re-ID

Cross-domain

4

# Sketch-photo person re-identification

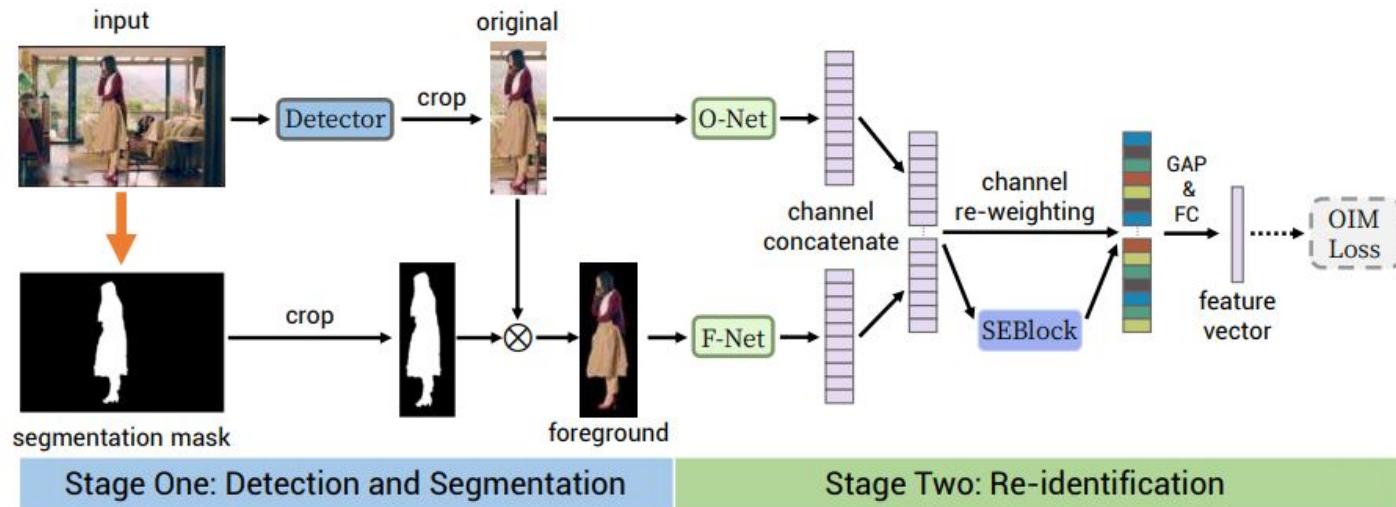


## Contributions

- A deep adversarial learning architecture to jointly learn identity features and domain invariable features.
- Filtering low-level features and remaining high-level semantic features.
- A sketch Re-ID dataset containing 200 persons, in which each person has one sketch and two photos.

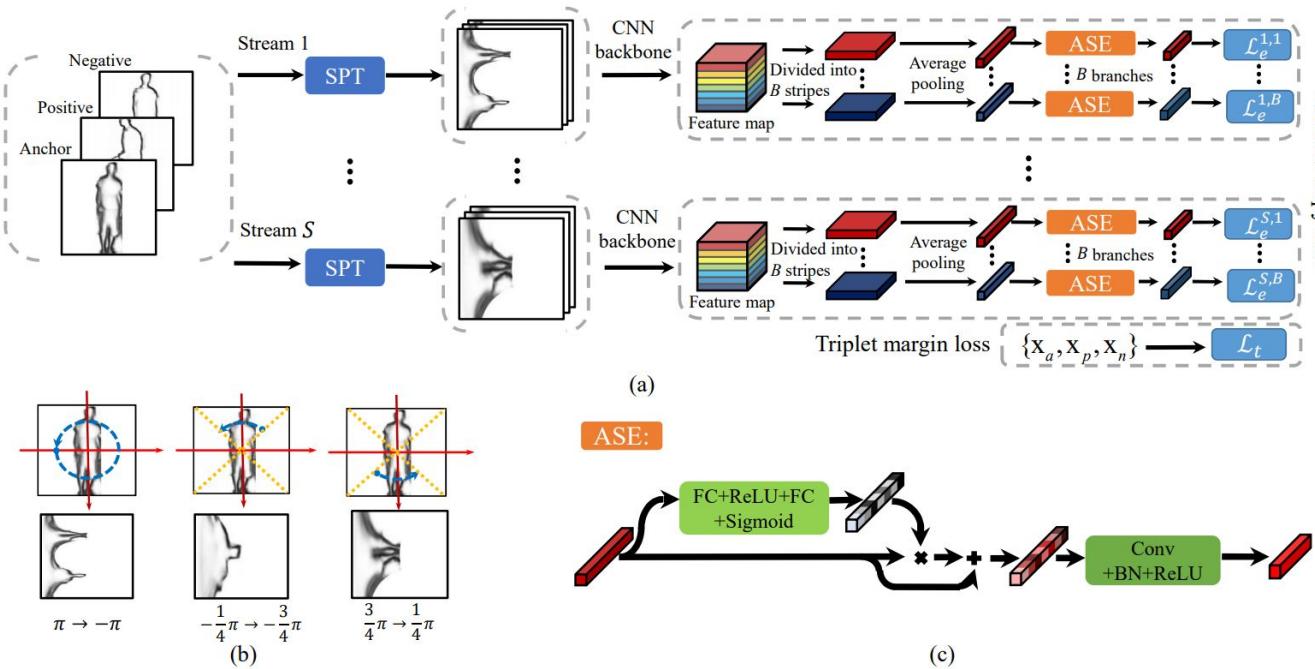
# Pedestrian detection + Person Re-identification

- Person search; person re-identification; person detection
- Contradict between detection and re-ID
- Background information



# Code (only dataset?)

## Person re-identification with clothing change



### Contributions:

- Find **contour sketches** are much more effective for person re-id under **moderate clothing change**.
- Introduce **SPT** to select the relatively invariant and discriminative contour patterns and **ASE** to explore angle-specific fine-grained discriminant features.
- PRCC: A new person re-id dataset with moderate clothing changes. (221 persons)

# Code (matlab) Bayesian based

Database



Query



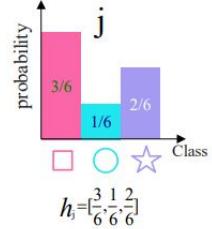
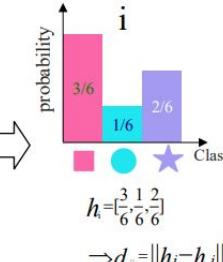
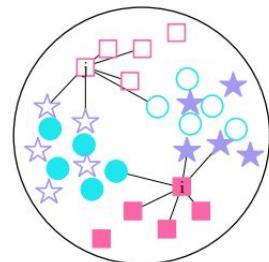
Returned



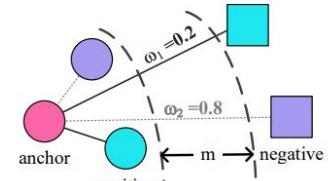
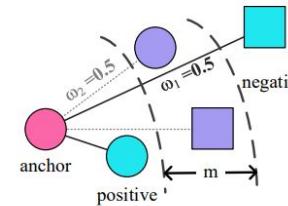
Single-domain retrieval



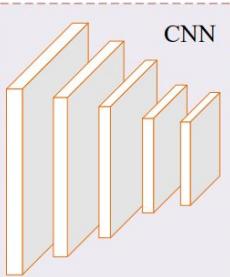
Cross-domain retrieval



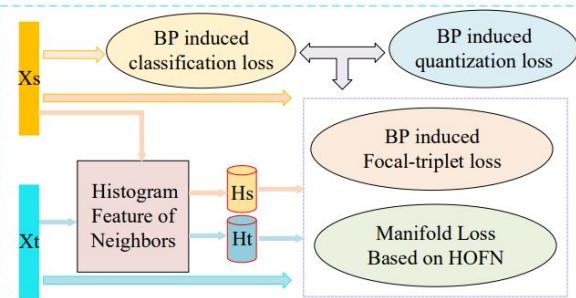
$$\Rightarrow d_{ij} = \|h_i - h_j\|^2 = 0$$



Source domain images

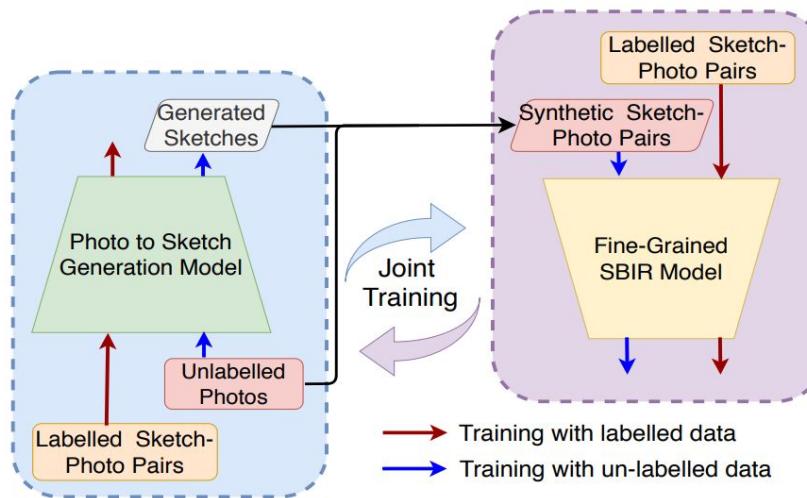


Feature extraction



Compactness

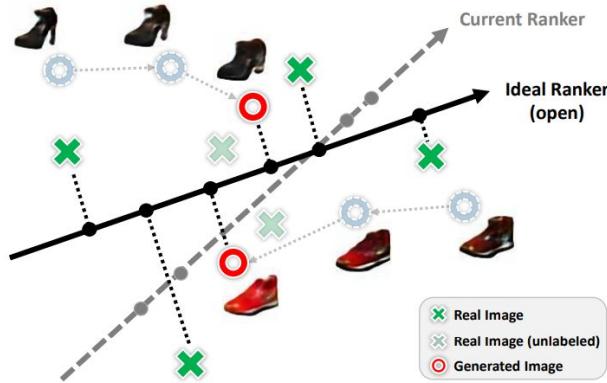
# More photos are all you need



# Relative Attributes

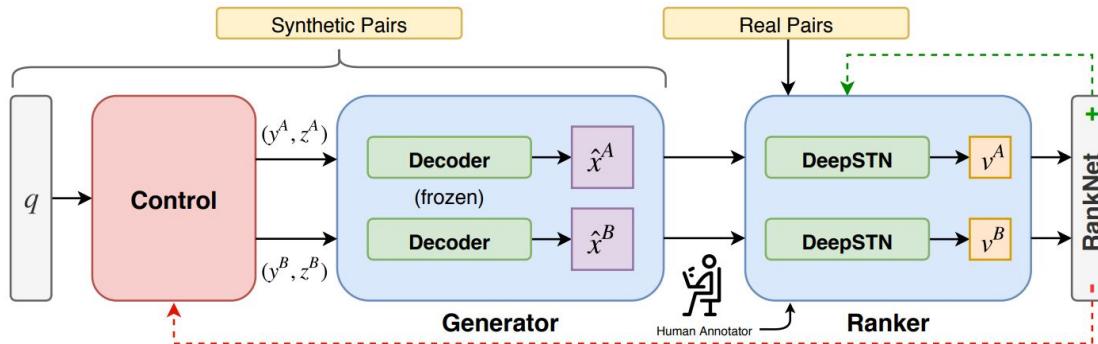
1

# ATTIC



- Contributions

- A new approach to active learning in which pairs of samples are newly created so as to best **augment a ranker's training data**. Active generation focuses attention on novel training images that rapidly improve generalization - even after all available real images and their labels are exhausted.

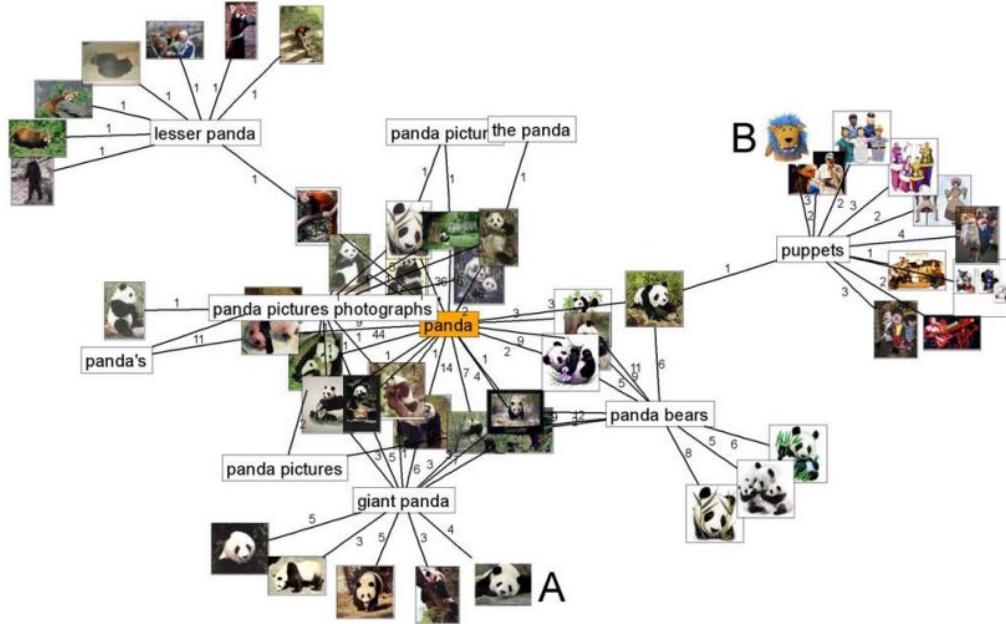


# Click based search

2

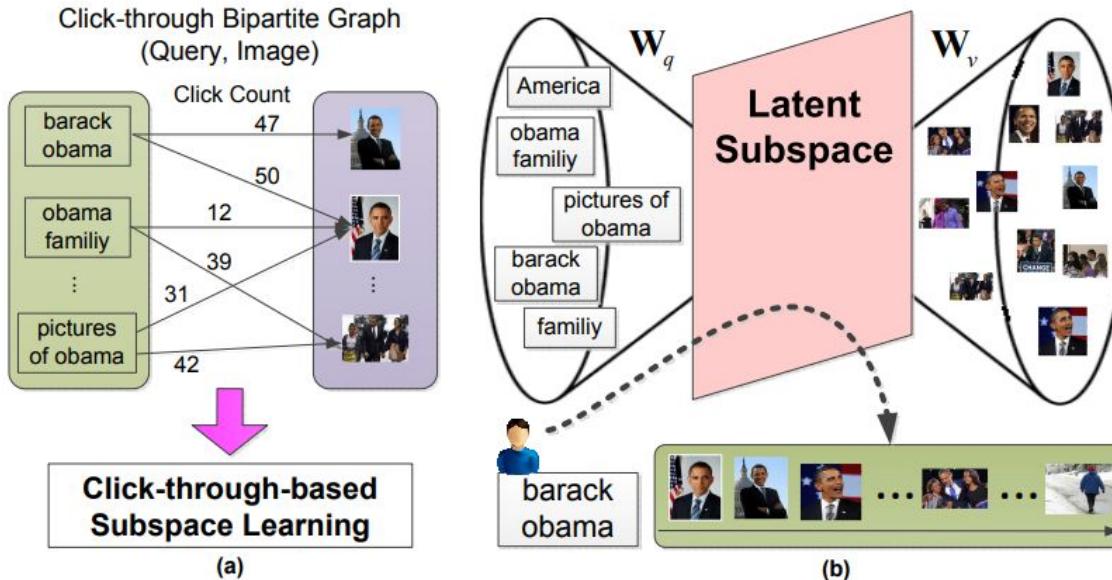
# Random walks on Click Graph

- Sparsity problem (documents that are relevant but not clicked) in clicks data



# Common latent subspace (text to image)+ user click data

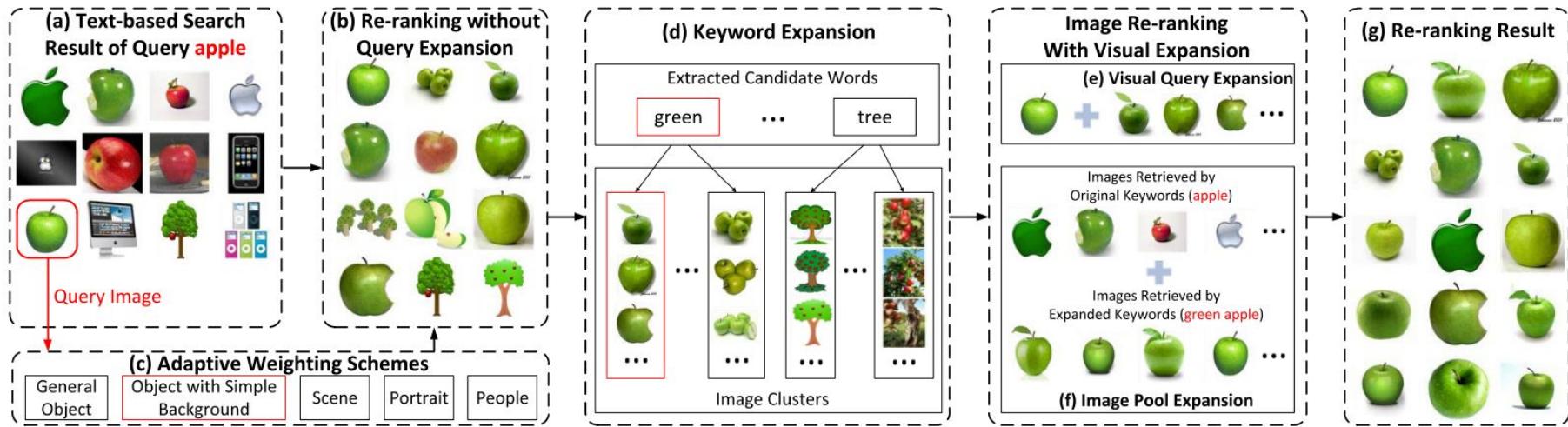
- Try to solve the ‘semantic gap’



[Pan, et al... Click-through-based Cross-view Learning for Image Search , SIGIR 2014 \(87\)](#)

[Pan, et al... Click-through-based Subspace Learning for Image Search , MM 2014 \(poster\) \(15\)](#)

# One click expanse



# Video Retrieval

1

# Video Moment Retrieval

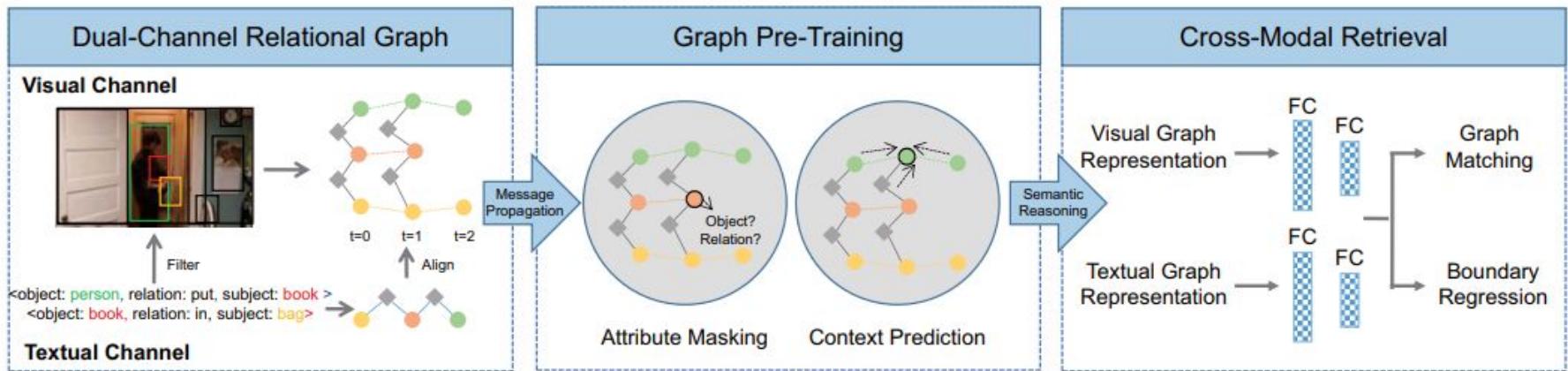
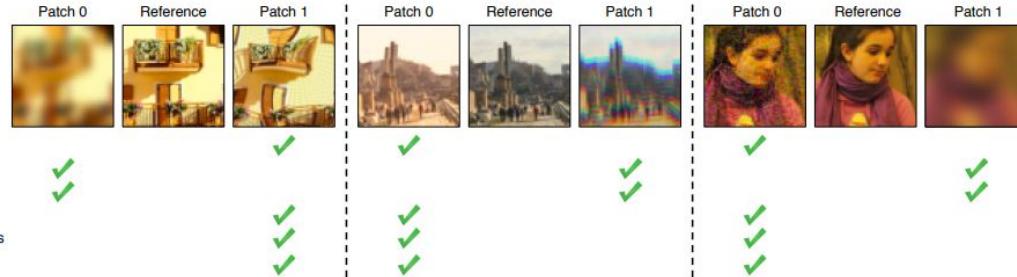
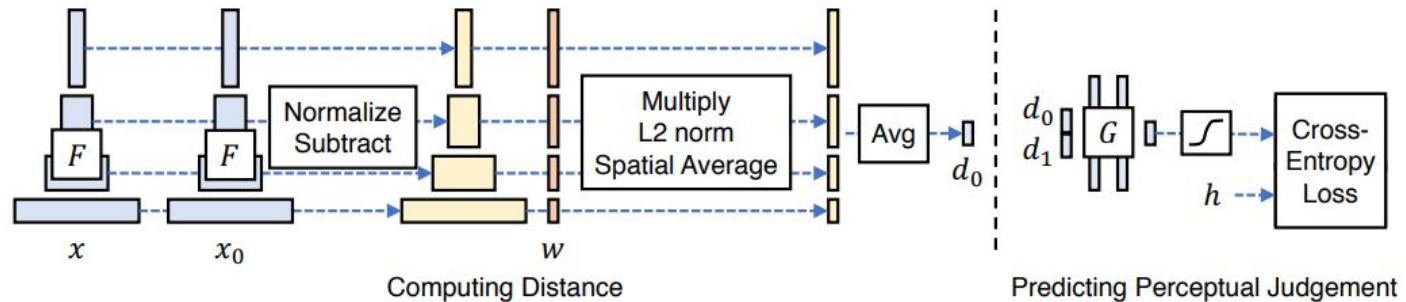


Figure 2: The graphical representation of our proposed MMRG framework. The input is an untrimmed video and its query sentence, while the output is the alignment score and location offset.

# Human perceptual Metric

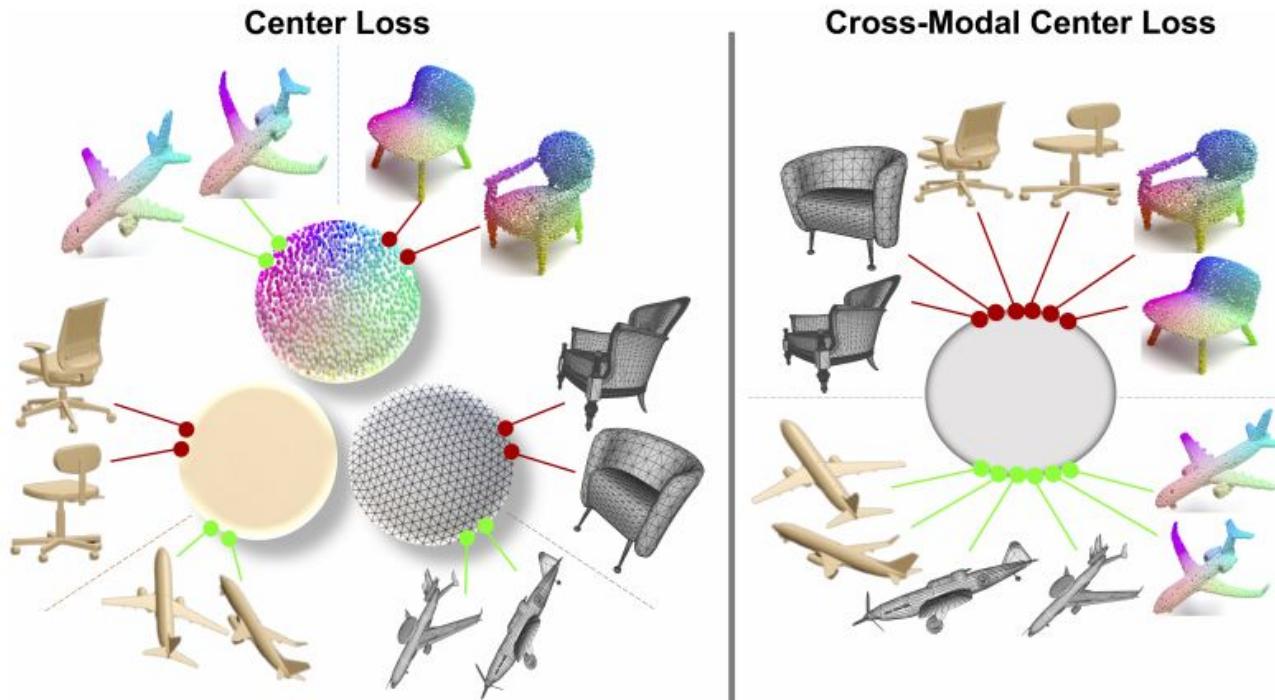
1

# DataSet + Much experiment

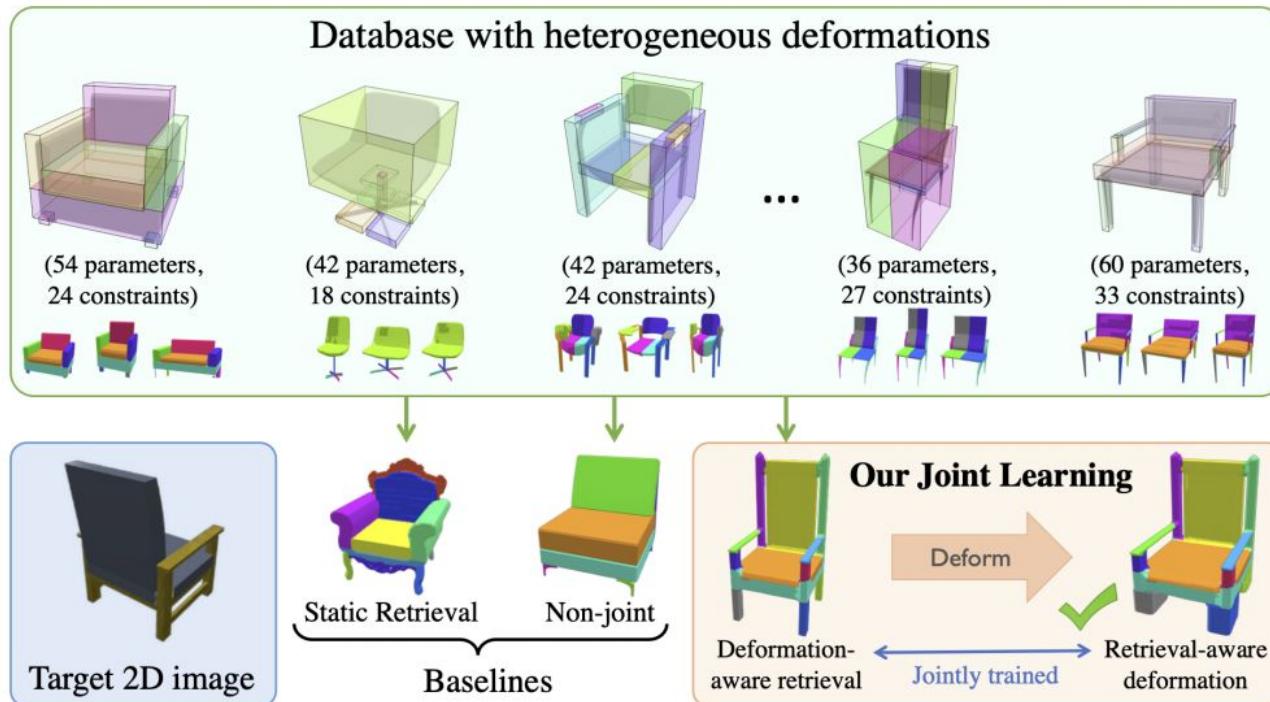


Special modal;  
Multi-Modal; 3D

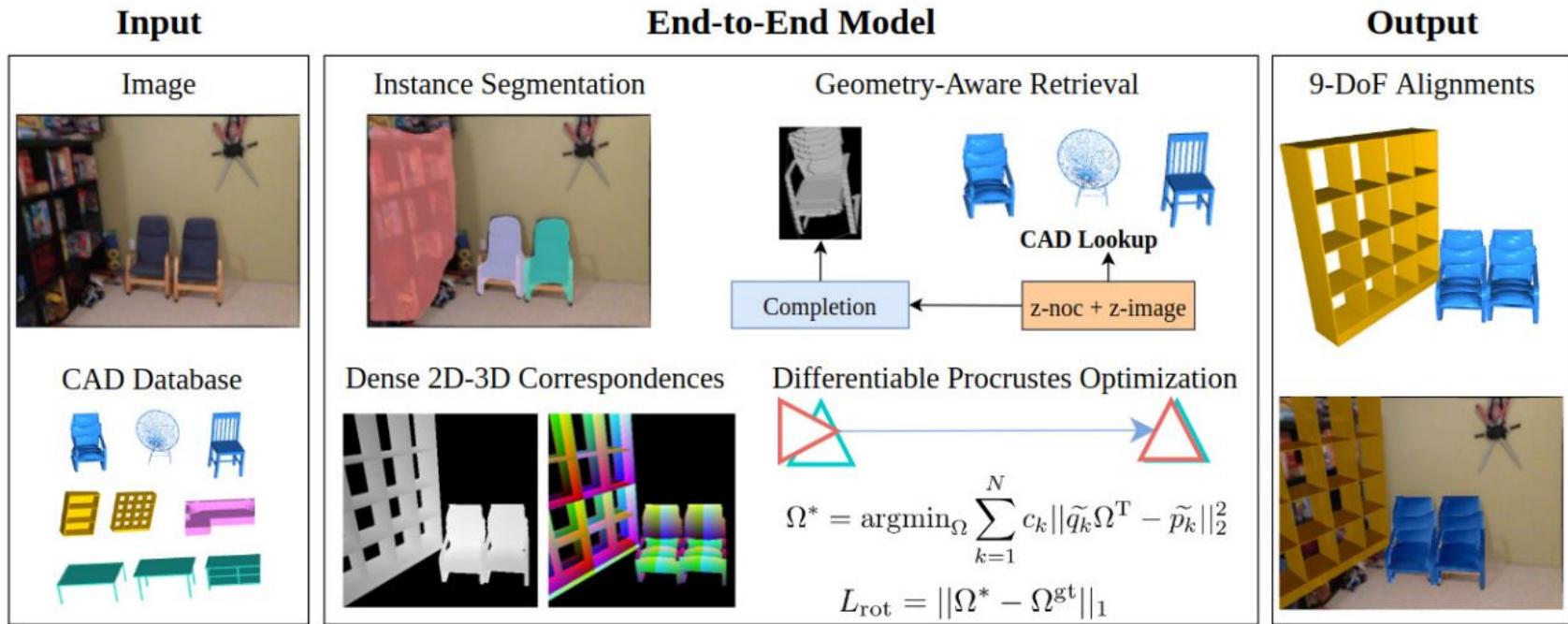
# Cross-Modal Retrieval



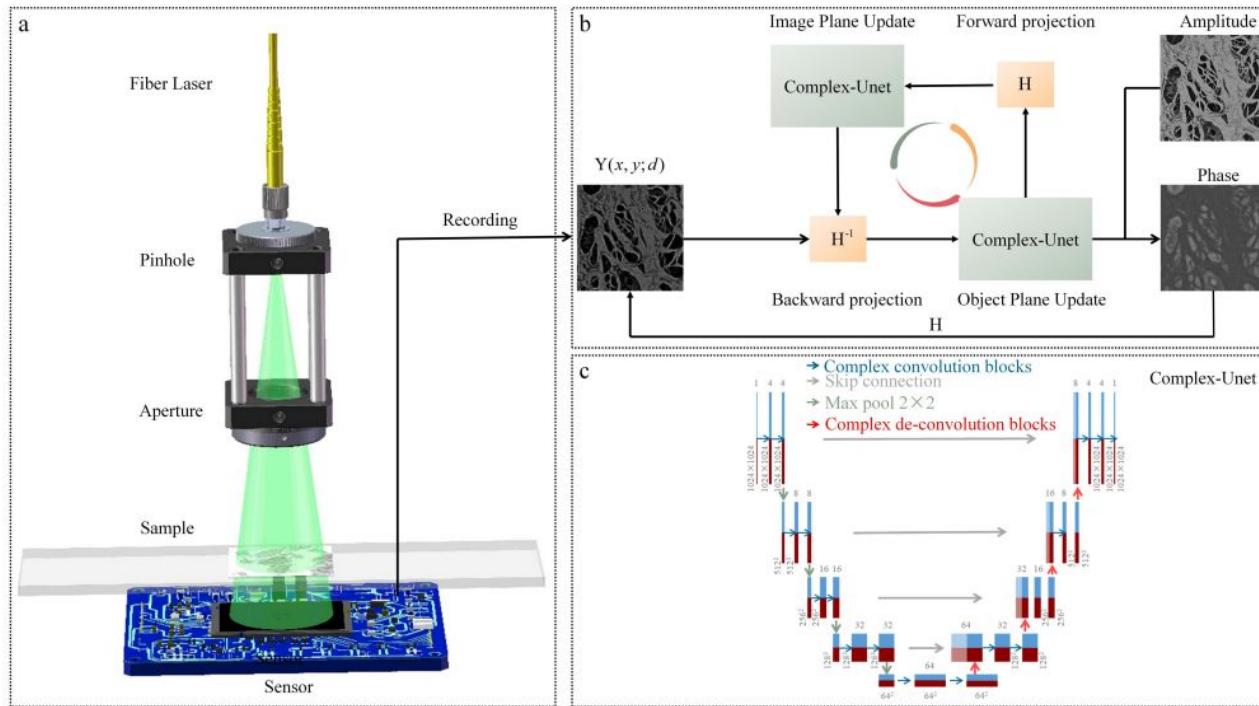
# 3D shape Retrieval



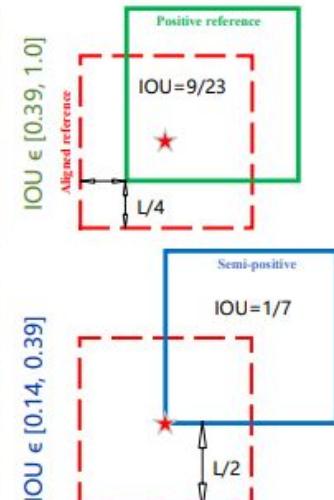
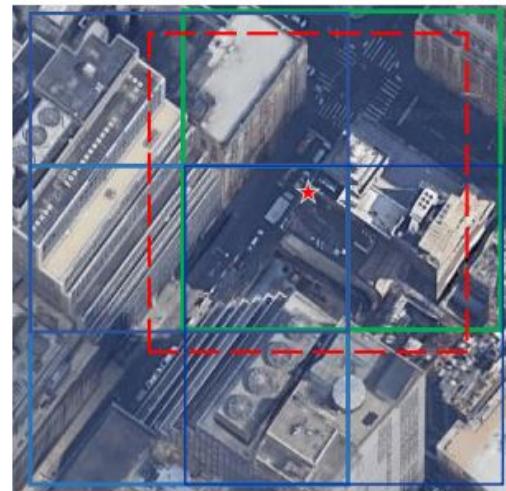
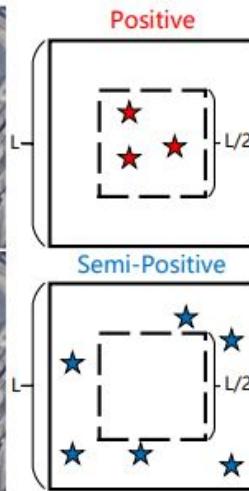
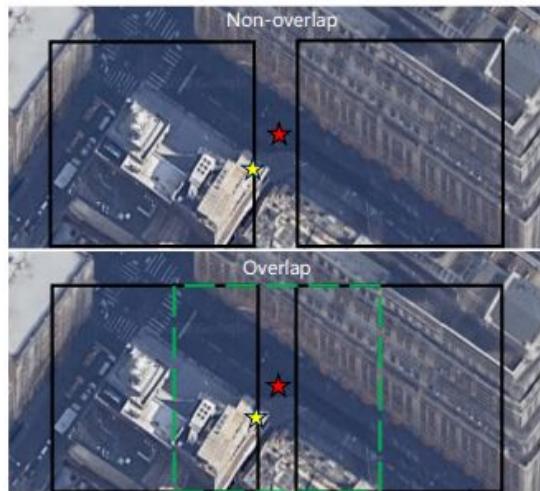
# CAD Retrieval



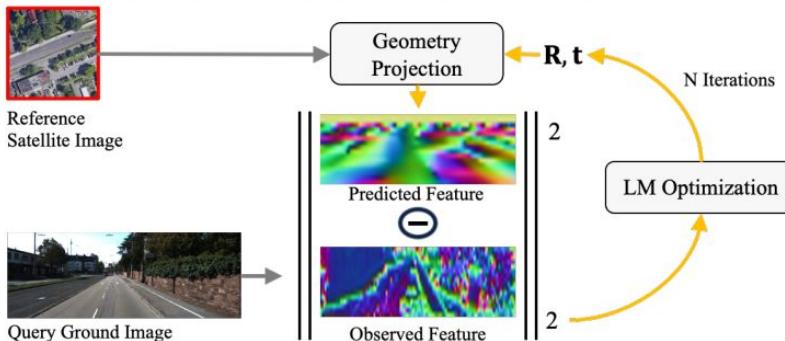
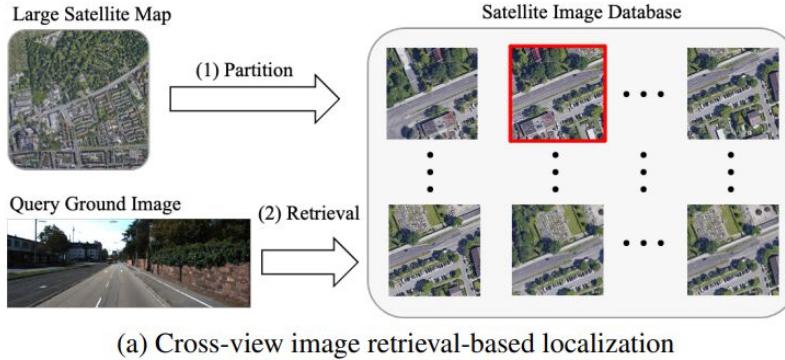
# Physics Phase Retrieval



# Geo-localization



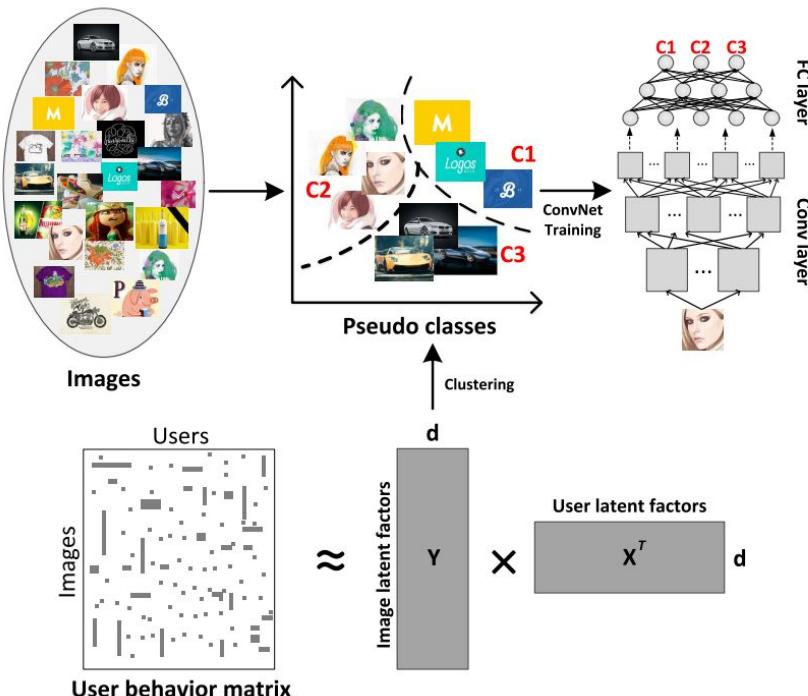
# Vehicle Localization; Satellite Image



# Interactive Image Search

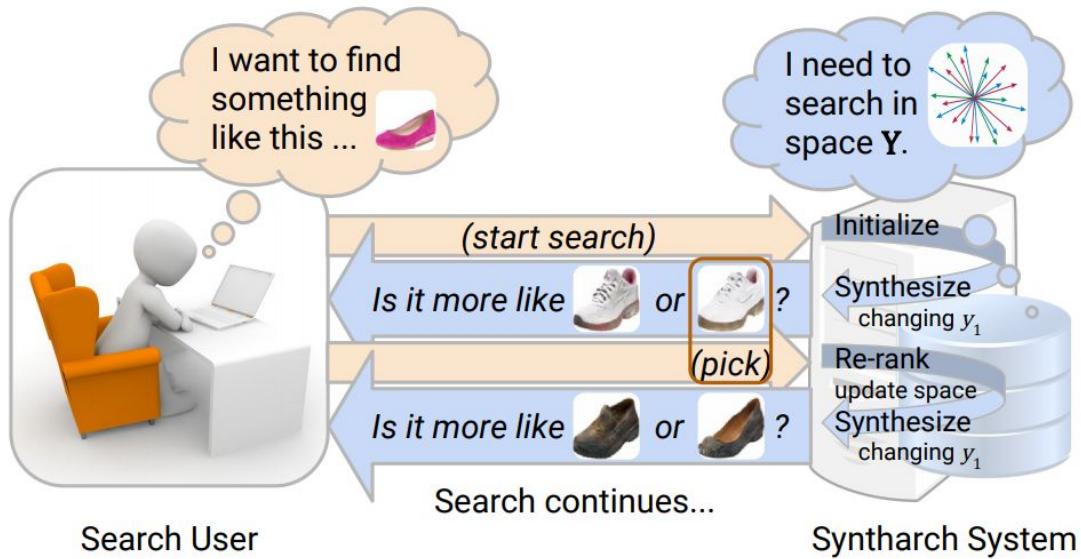
3

# Image recommendation



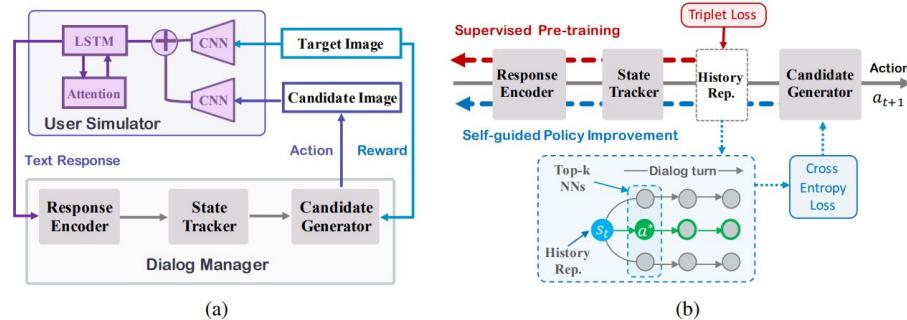
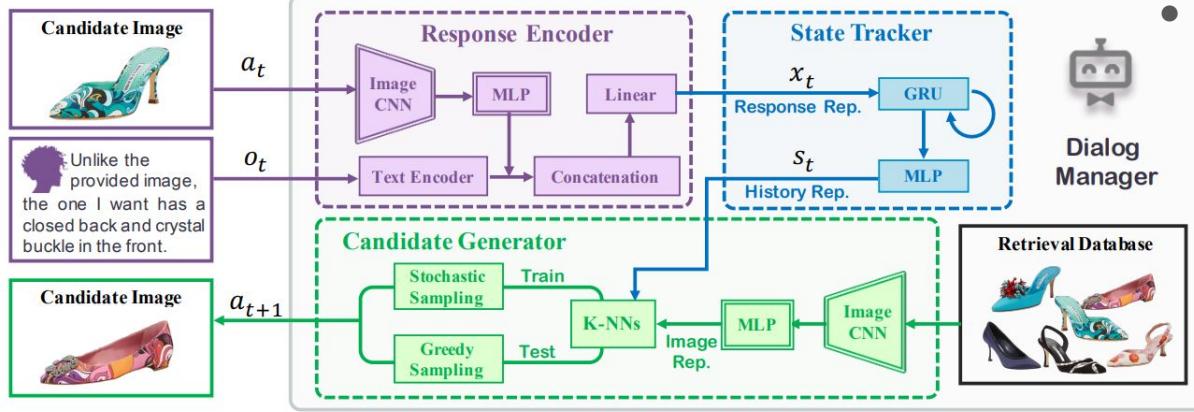
- Contributions

- Completely forgo the use of category labels in existing feature learning paradigms. Instead, they use user behavior data collected on social media.
- Test on large-scale data collected from a real-world social media website.



# Code

# Dialog-based Interactive Image Retrieval



## Contributions

- A new vision/NLP task.
- Dialog manager architecture with triplet loss and model-based policy improvement
- A new computer vision task, relative image captioning.
- A new dataset.





# Dateset

6

# Code Oxf5k

- Detail

- 5062 images collected from **Flickr** by searching for particular **Oxford landmarks**.
- manually annotated; 11 different landmarks, each represented by **5 possible queries**.



55 queries

Search similar images from  
5062 images pool



# Code

## Paris and Flickr1

Dataset	Number of images	Number of features
<i>Oxford</i>	5,062	16,334,970
<i>Paris</i>	6,300	20,219,488
<i>Flickr1</i>	99,782	277,770,833
Total	111,144	314,325,291

Table 1. The number of descriptors for each dataset.

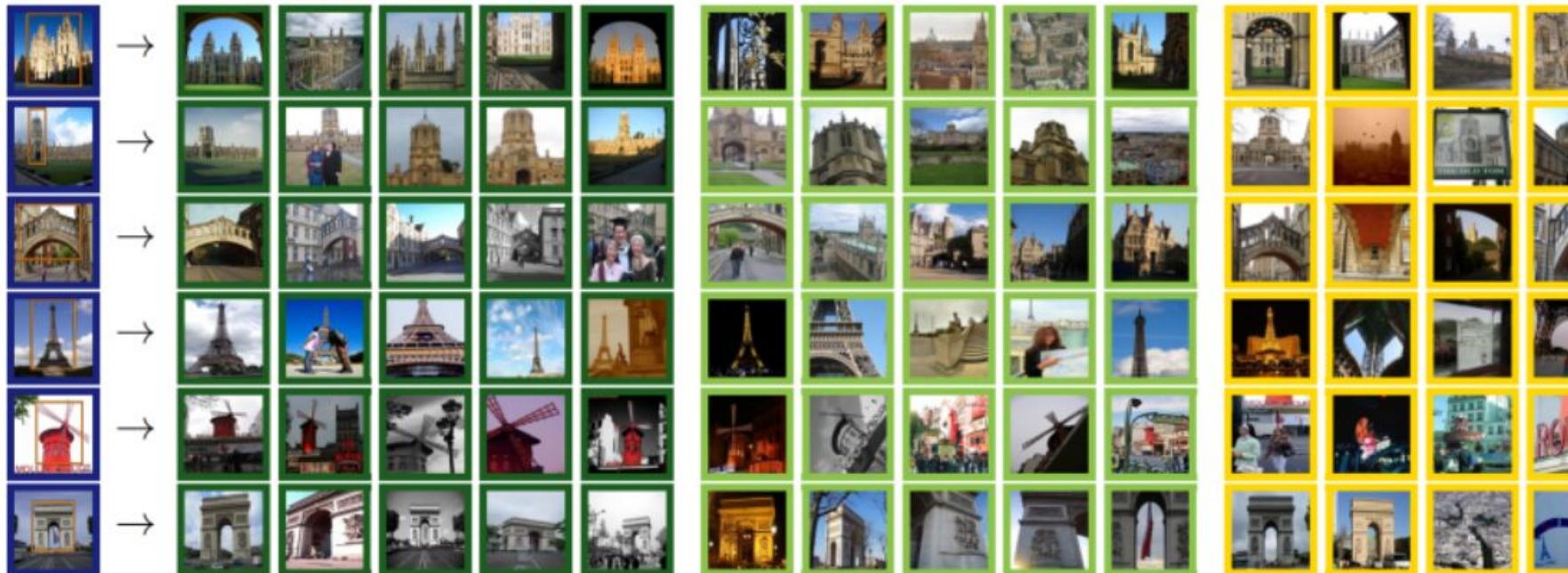


# Code

## R-Oxf & R-Par

- Detail

- ROxford: 4993 images + 70 queries.



# Code GLDv2

- Detail

- **GLDv2-retrieval:** 1129 queries(379 validation and 750 testing) and 762k database images
- **GLDv2-recognition:** 118K test (41k validation and 77k testing) and 4M training images from 203k landmarks.

Dataset name	Year	# Landmarks	# Test images	# Train images	# Index images	Annotation collection	Coverage	Stable
Oxford [41]	2007	11	55	-	5k	Manual	City	Y
Paris [42]	2008	11	55	-	6k	Manual	City	Y
Holidays [28]	2008	500	500	-	1.5k	Manual	Worldwide	Y
European Cities 50k [5]	2010	20	100	-	50k	Manual	Continent	Y
Geotagged StreetView [32]	2010	-	200	-	17k	StreetView	City	Y
Rome 16k [1]	2010	69	1k	-	15k	GeoTag + SfM	City	Y
San Francisco [14]	2011	-	80	-	<b>1.7M</b>	StreetView	City	Y
Landmarks-PointCloud [35]	2012	1k	10k	-	205k	Flickr label + SfM	Worldwide	Y
24/7 Tokyo [55]	2015	125	315	-	1k	Smartphone + Manual	City	Y
Paris500k [60]	2015	13k	3k	-	501k	Manual	City	N
Landmark URLs [7, 22]	2016	586	-	140k	-	Text query + Feature matching	Worldwide	N
Flickr-SfM [44]	2016	713	-	120k	-	Text query + SfM	Worldwide	Y
Google Landmarks [39]	2017	30k	<b>118k</b>	1.2M	1.1M	GPS + semi-automatic	Worldwide	N
Revisited Oxford [43]	2018	11	70	-	5k + 1M	Manual + semi-automatic	Worldwide	Y
Revisited Paris [43]	2018	11	70	-	6k + 1M	Manual + semi-automatic	Worldwide	Y
Google Landmarks Dataset v2	2019	<b>200k</b>	<b>118k</b>	<b>4.1M</b>	762k	Crowdsourced + semi-automatic	Worldwide	Y

Table 1: Comparison of our dataset against existing landmark recognition/retrieval datasets. “Stable” denotes if the dataset can be retained indefinitely. Our Google Landmarks Dataset v2 is larger than all existing datasets in terms of total number of images and landmarks, besides being stable.

# Code

## Attribute Discovery Dataset

- Detail

- 4 broad shopping categories (**bags, earrings, ties, and shoes**).
- 37795 images: 9145 handbags; 14765 shoes; 9235 earrings; 4650 ties.
- Collected from like.com
- images and associated textual descriptions



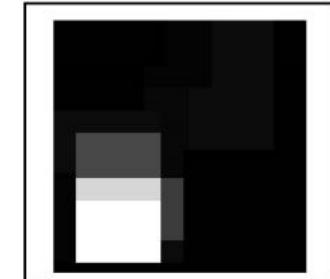
Dazzle after dark with Judith Leiber's decadent oversized crystal-embellished silver-tone clutch. Carry this fabulous extra to add high-octane glamour to an LBD and teetering heels. Shown here with an Emilio Pucci dress and Givenchy shoes.



The 12K pink and green gold leaves gently cascade down on these delicate beaded 10K gold earrings.



Rock and roll in the Signature. The sm leather upper with atop a 1 inch plat ankle with a gold l heel. Sizzle in the:



# Code

## Relative captioning dataset

- Detail
  - Shoes dataset. 14658 images
  - Relative\_captions\_shoes.json 10751 captions, with one caption per pair of images.
  - Captions\_shoes.json: captions on 3600 shoe images



(b)

# UT-Zap50K

[Yu, et al... Fine-grained visual comparisons with local learning, CVPR 2014 \(247\)](#)

[Yu, et al... Semantic jitter: Dense supervision for visual comparisons via synthetic images, ICCV 2017 \(56\)](#)

# DeepFashion



Figure 2. Example images of different categories and attributes in DeepFashion. The attributes form five groups: texture, fabric, shape, part, and style.

# DARN



(a) Query Image



(b) Top-6 Retrieval Results

To read