

MIT
COMPUTER
VISION



6.819 / 6.869: Advances in Computer Vision

Image Retrieval:

Retrieval: Information, images, objects, large-scale

Website:

<http://6.869.csail.mit.edu/fa15/>

Instructor: Yusuf Aytar

Lecture TR 9:30AM – 11:00AM
(Room 34-101)

What is Image Retrieval ?

User



Query

Text

Fall in Boston

Image



Speech



Retrieval Results



Applications

Art Retrieval



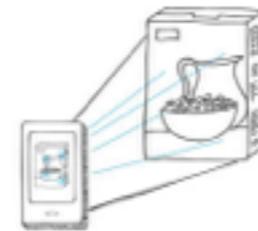
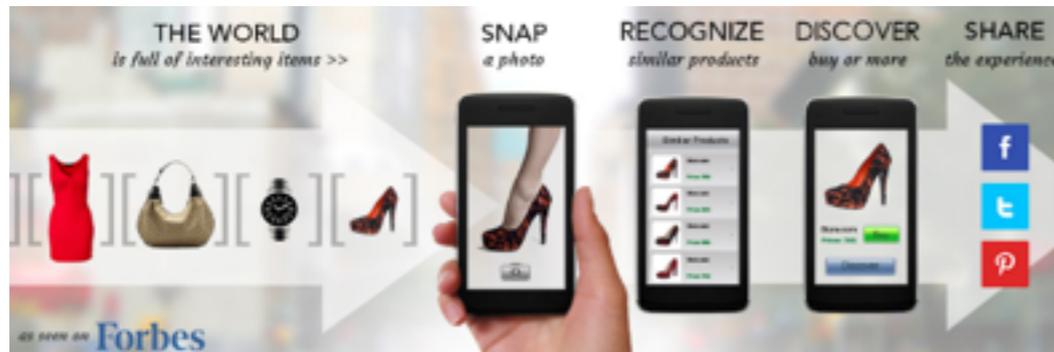
Medical Image Retrieval



Product Image Retrieval
(Reviews, other prices etc.)



Fashion Image Retrieval



Overview

Information Retrieval

Bag of Words, TF-IDF, Cosine Similarity, Inverted Index

Object Instance Retrieval

Bag of Visual Words, Video Google, Object Instance Retrieval

Fast Object Detection/Retrieval

Fast detection, Part representations, Generalization from exemplar

Large Scale Image Search

KD-trees, Locality Sensitive Hashing, Semantic Hashing, Compact Codes

Information Retrieval

Bag of Words (BOW)

The last duel

After quarrelling over a bank loan, two men took part in the last fatal duel staged on Scottish soil. BBC News's James Landale retraces the steps of his ancestor, who made that final challenge.

Doc-1

Bank	: 1
Loan	: 1
Water	: 0
Farmer	: 0

West Bank water row

Palestinians have accused Israel of diverting water away from their towns in order to keep Jewish settlements in the occupied territories fully supplied. Israel denies the charge saying Palestinian farmers are to blame for using illegal connections to irrigate their fields.

Doc-2

Bank	: 1
Loan	: 0
Water	: 2
Farmer	: 0

A widely used document representation method

Term Frequency (TF)

The last duel

Uncovering the trenches
I know quite a lot about these people, through documentary evidence relating to them and contemporary accounts of their times. I can identify the ship my ancestor served on in the Seven Years' War, his father's house on Holy Island and probably the...

Doc-3

Bank	: 0
Loan	: 1
Water	: 0
Farmer	: 0

West Bank water row

Morning sickness
Just a few of the opening bars are enough to transport many unsuspecting souls back to the school assembly halls of their childhood...

Doc-4

Bank	:0
Loan	:3
Water	:1
Farmer	:1

	Documents			
Lexicon	Doc-1	Doc-2	Doc-3	Doc-4
Bank	1	1	0	0
Loan	1	0	0	3
Water	0	2	1	1
Farmer	0	0	0	1

Normalization



	Documents			
Lexicon	Doc-1	Doc-2	Doc-3	Doc-4
Bank	0.5	0.33	0	0
Loan	0.5	0	0	0.6
Water	0	0.66	1	0.2
Farmer	0	0	0	0.2

Inverse Document Frequency (IDF)

IDF of i^{th} word: $idf_i = \log\left(\frac{n}{df_i}\right)$

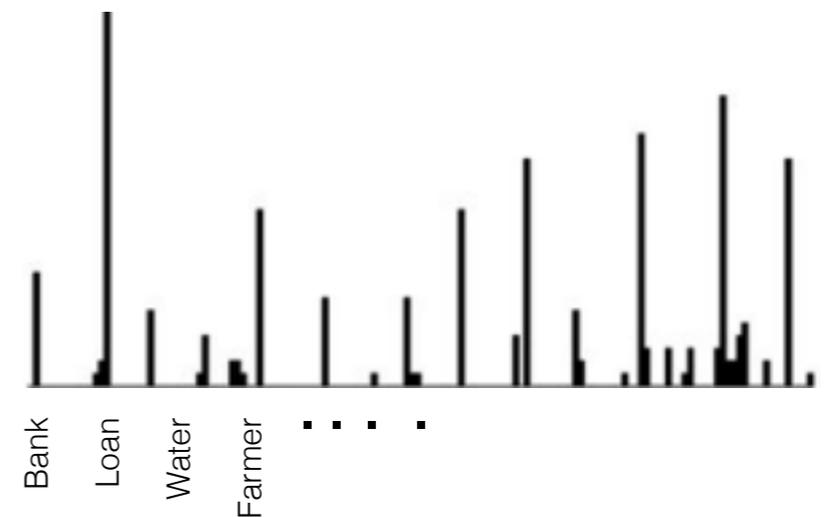
The last duel

After quarrelling over a bank loan, two men took part in the last fatal duel staged on Scottish soil. BBC News's James Landale retraces the steps of his ancestor, who made that final challenge.

Doc-1



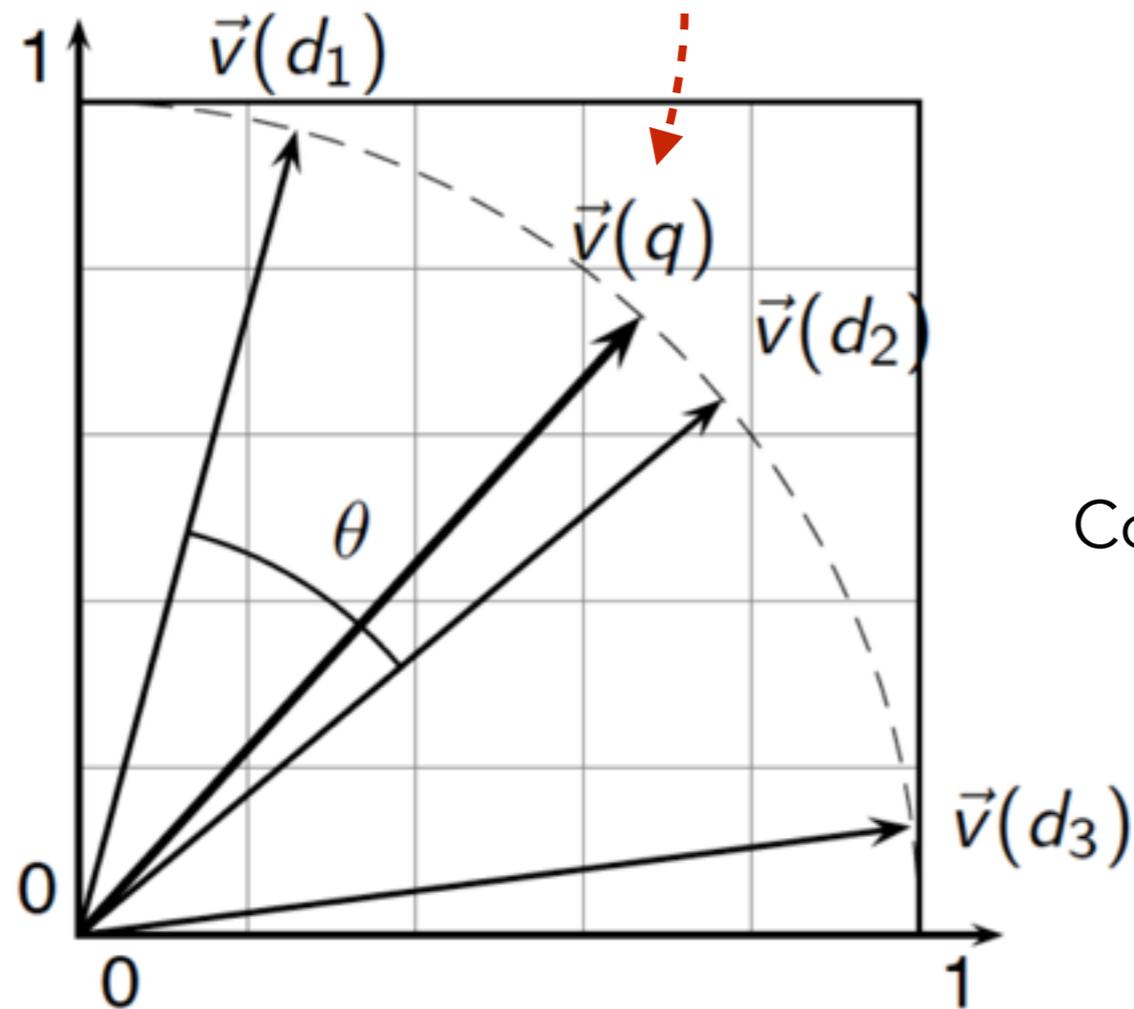
$$\vec{v}(d_1) =$$



$tf \times idf$

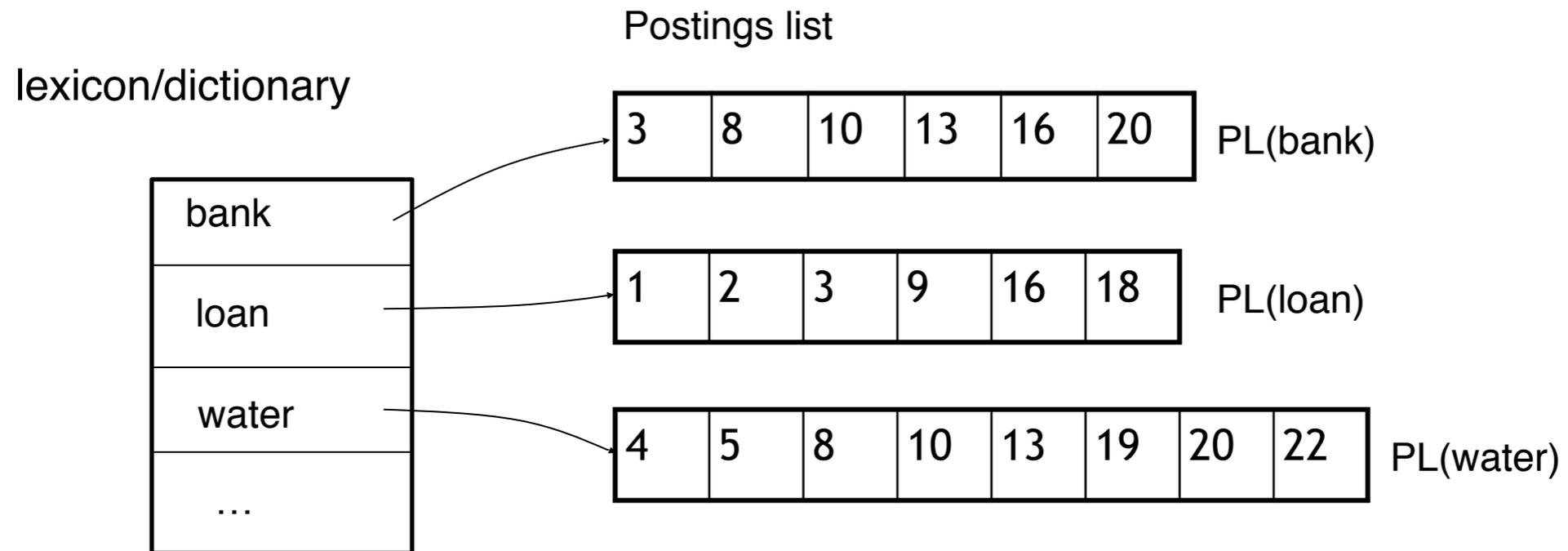
Cosine Similarity

Query: **fall in Boston**



$$\text{Cosine Similarity Score} = \vec{v}(q)^\top \vec{v}(d_1)$$

Inverted Index



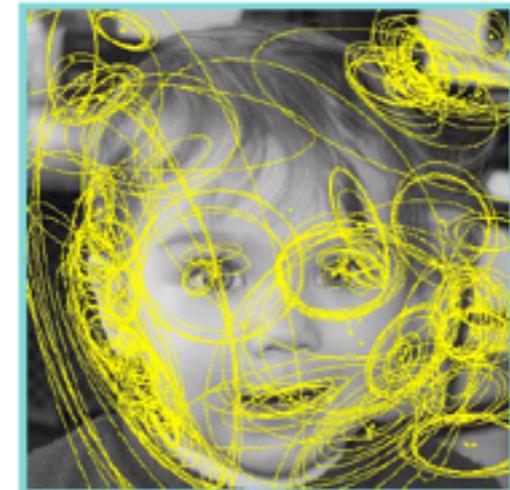
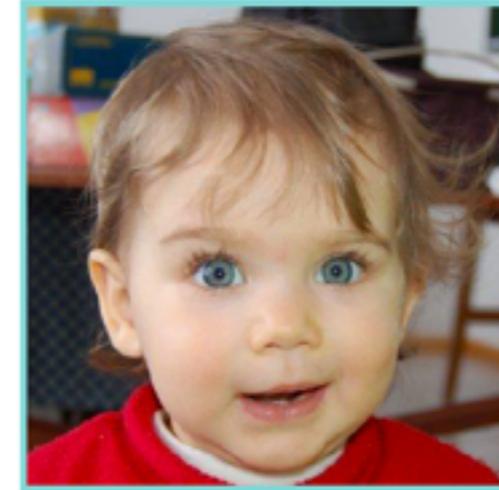
Allows quick lookup of document ids with a particular word

Bag of Words & Object Instance Retrieval

Feature Detectors

Roles of the detector:

- provide invariance to transformations
- **reduce the number of descriptors**



Popular detectors:

- **Maximally Stable Extremal Regions (MSER)**
Matas, Chum, Urban, Pajdla, "Robust wide-baseline stereo from maximally stable extremal regions", BMVC'02.
- **Difference of Gaussians (DoG)**
Lowe, "Distinctive image features from scale-invariant keypoints", IJCV'04.
- **Harris-Affine and Hessian-Affine**
Mikolajczyk, Schmid, "Scale and affine invariant interest point detectors", IJCV'04.
→ See also Mikolajczyk et al., "A comparison of affine region detectors", IJCV'05.

Dense descriptors are also possible

- **Mainly for classification → let the classifier decide**
Leung, Malik, "Representing and recognizing the visual appearance of materials using 3D textons", IJCV'01.
- **But also for image/scene/object retrieval**
Gordo, Rodriguez, Perronnin, Valveny, "Leveraging category-level labels for instance-level image retrieval", CVPR'12.

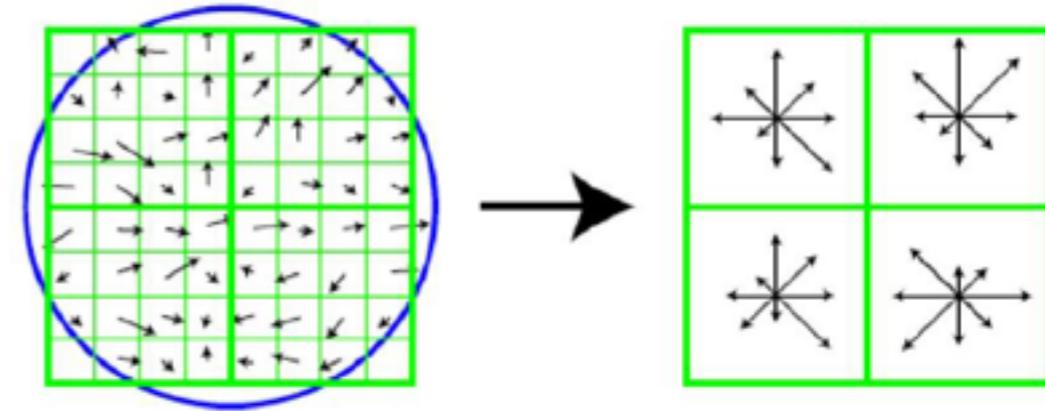
Feature Descriptors

Description of a patch after orientation/scale/photometric normalization

Most widely-used patch descriptor: SIFT

Lowe, "Distinctive image features from scale-invariant keypoints", IJCV'04.

- 8 orientations of the gradient → 128 dimensions
- 4x4 spatial grid



Many descriptors derive from SIFT:

- **More efficient: SURF**

Bay, Tuytelaars, Van Gool, "SURF: speeded up robust features", ECCV'06.

- **More compact: CHOG, DAISY**

Chandrasekhar et al, "Compressed histograms of gradients: a low-bit rate descriptor", IJCV'11.

Tola, Lepetit, Fua, "DAISY: an efficient dense descriptor applied to wide baseline stereo", TPAMI'10.

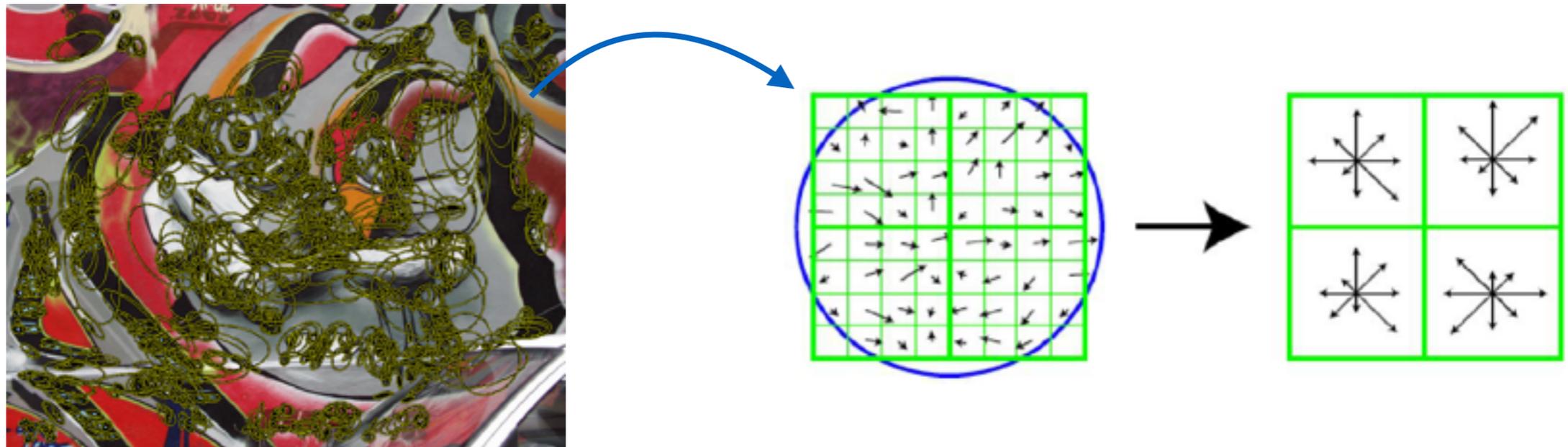
- **With color: color SIFT**

Van de Weijer, Schmid, "Coloring local feature extraction", ECCV'06.

Burghouts and Geseborek, "Performance evaluation of local colour invariants", CVIU'09.

Video Google

Feature Detectors / Descriptors

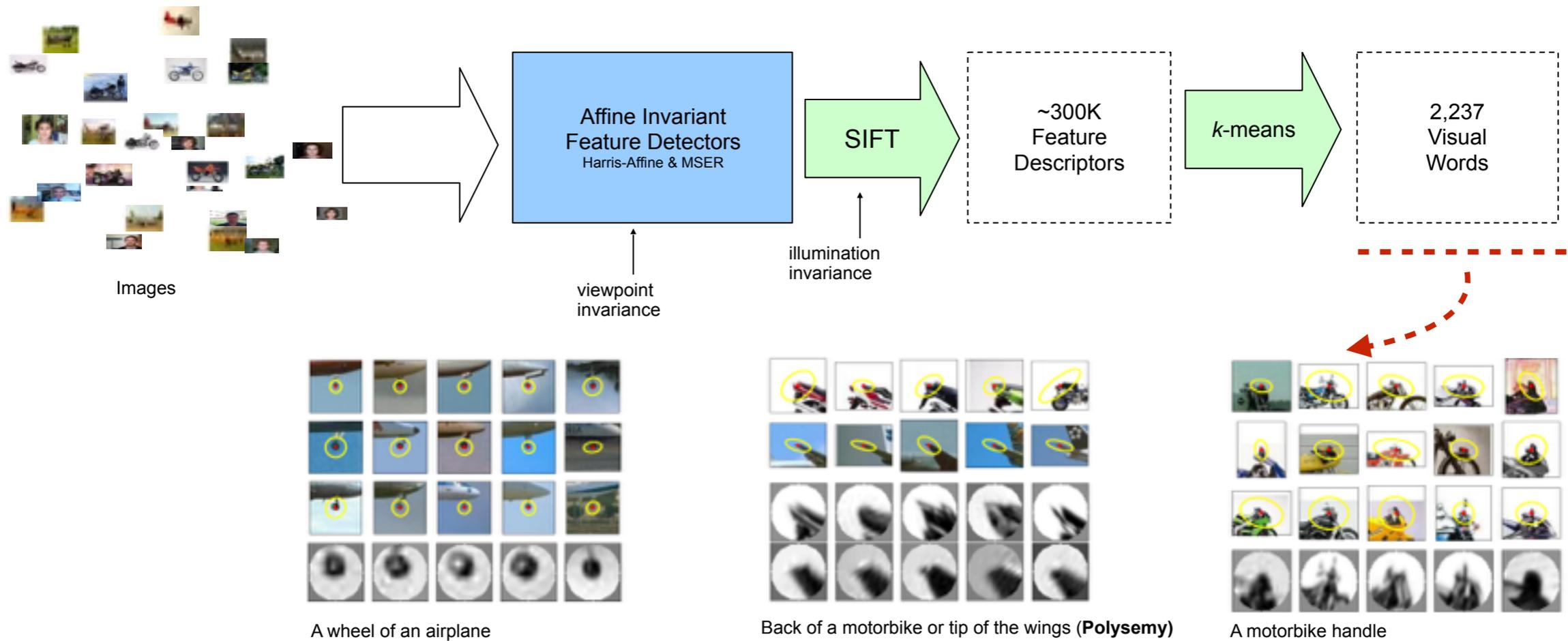


Harris-Affine & Hessian Affine as the **feature detectors**

SIFT as the **feature descriptor**

Video Google

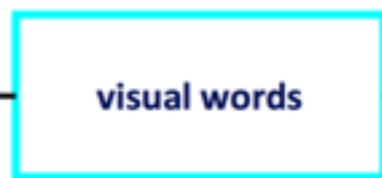
Bag of Visual Words



Set of SIFT descriptors



[Sivic03]



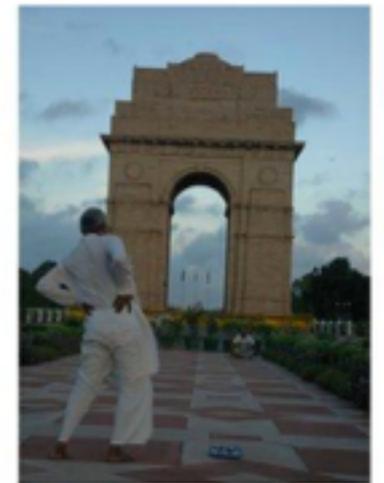
sparse frequency vector



Video Google

Large scale object instance retrieval

- Find all instances of the query object in a large scale dataset
- Do it **instantly** (< 1sec), and be robust to **scale, viewpoint, lighting, partial occlusion**



Video Google

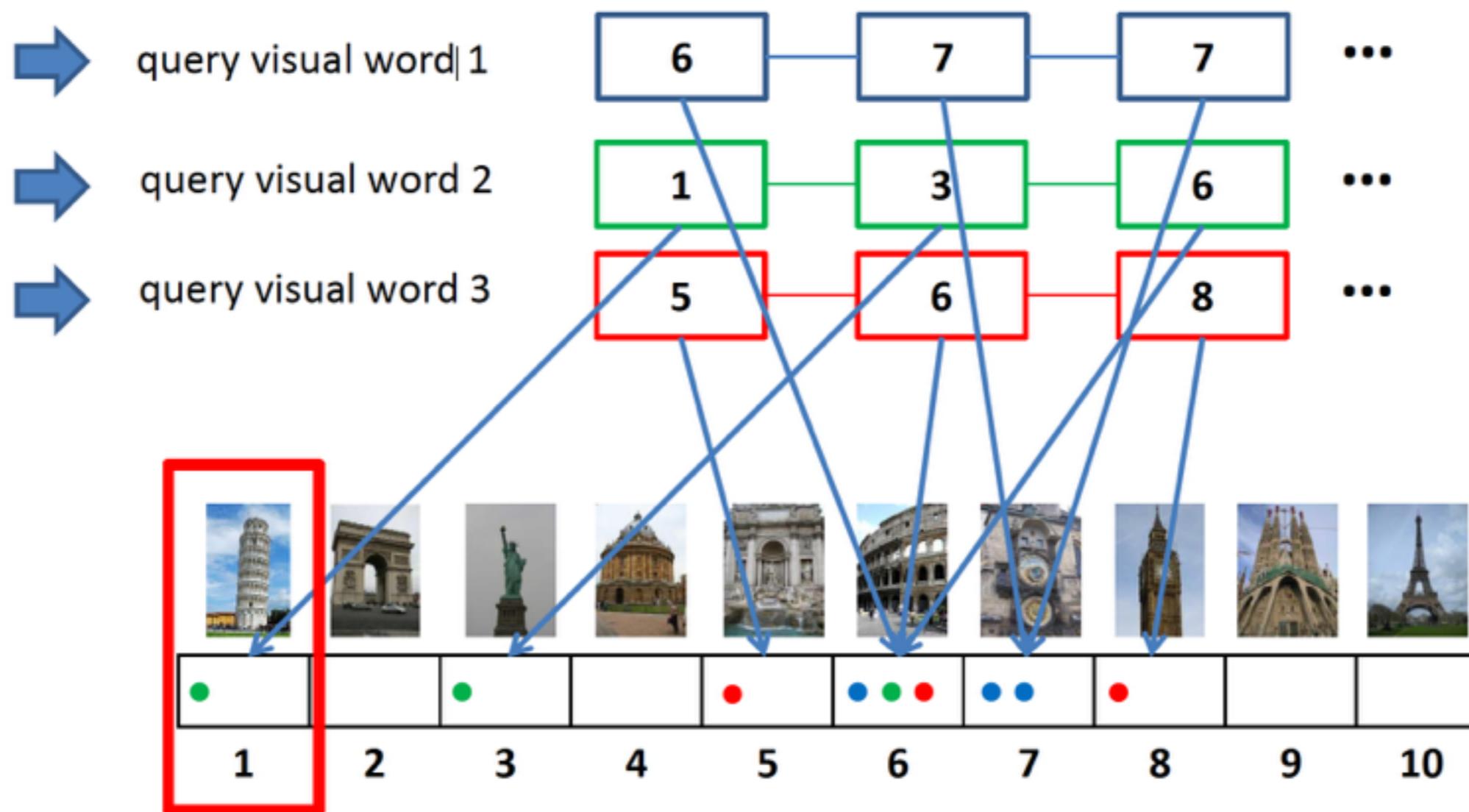
Particular object retrieval - Bag of visual words



Video Google

BOW + Inverted File Indexing

$$\text{score} = \frac{\mathbf{q}^T \mathbf{x}}{\|\mathbf{x}\|}$$



Spatial Verification

1. Image Query



2. Initial Retrieval Set



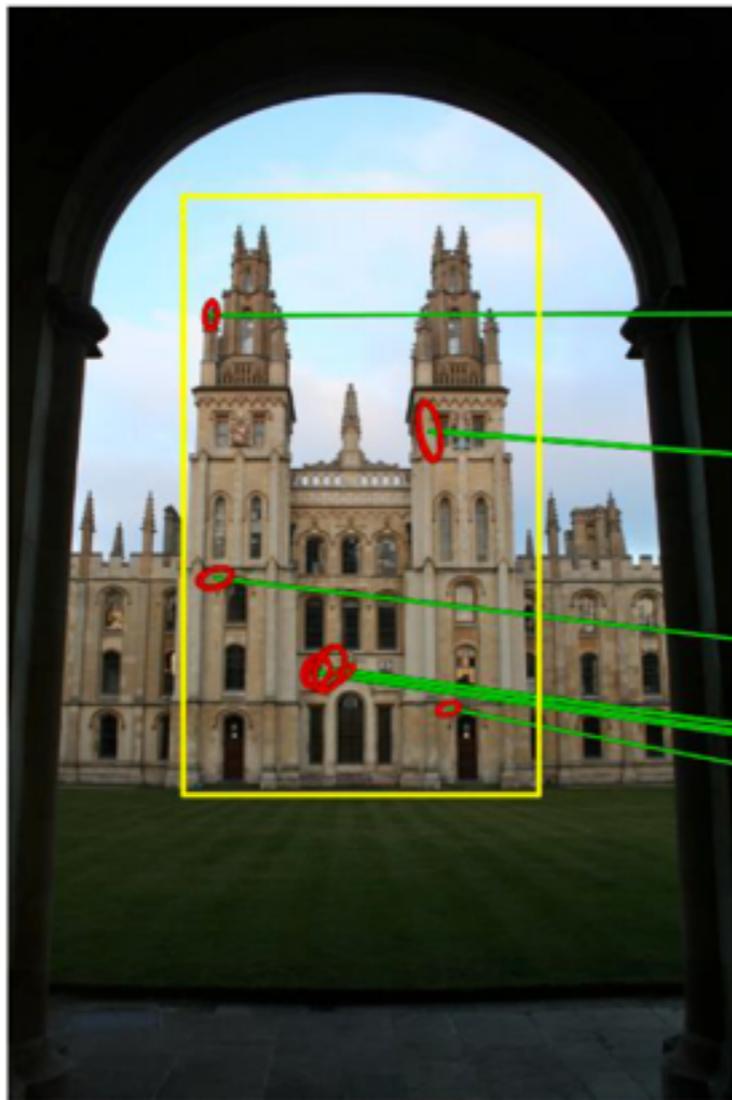
3. Spatial Verification



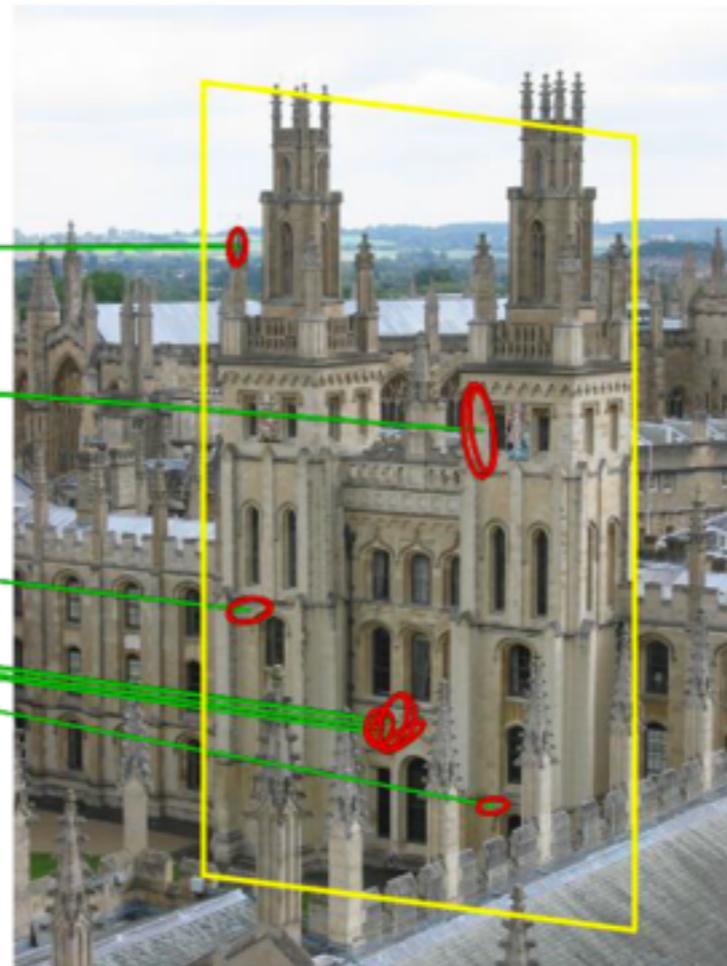
Query Expansion



Query Expansion



Query Image

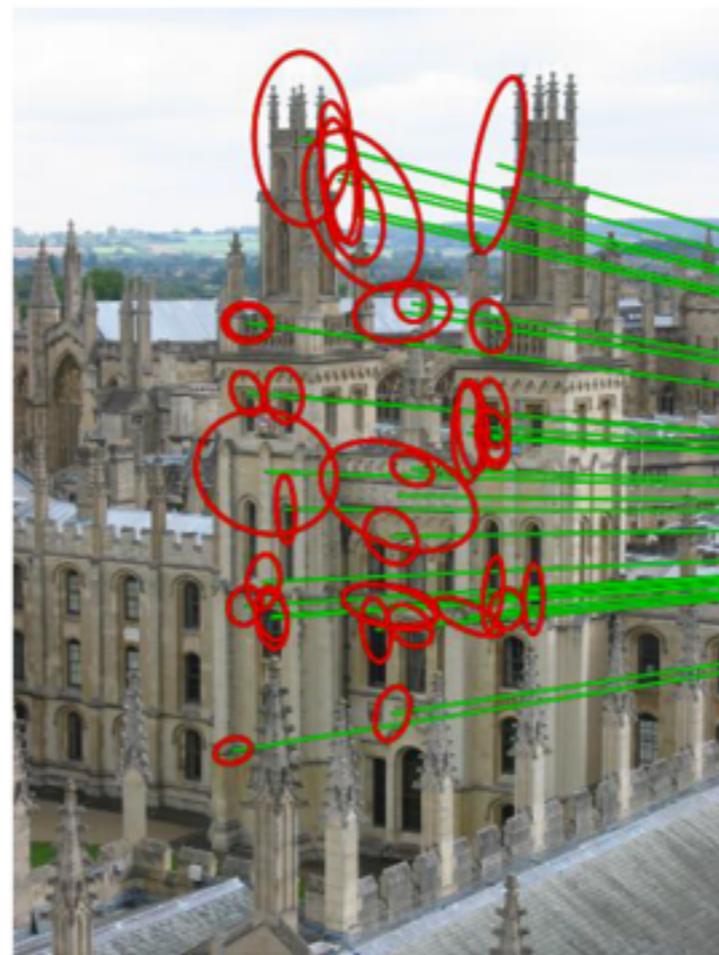
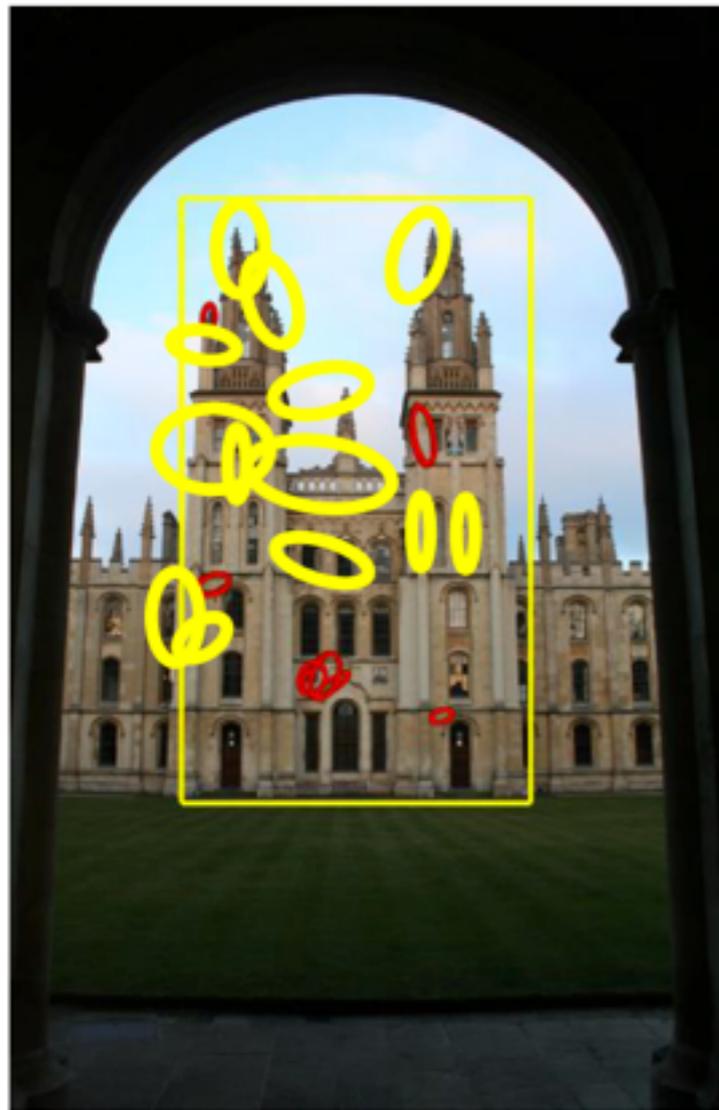


Retrieved image



Originally not retrieved

Query Expansion



Video Google - Object Instance Retrieval

videogoogle

Exploring Charade

Viewing frame 29300

Overview Explore shots
Prev Animate DivX Stream Thumbnails Search Next



Clear Search

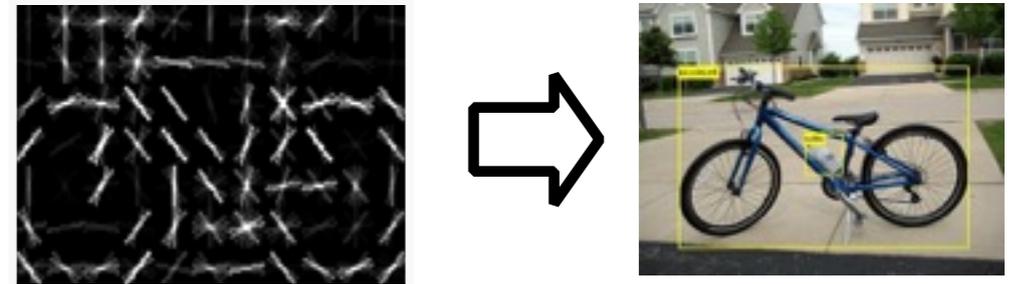


<http://www.robots.ox.ac.uk/~vgg/demo/>

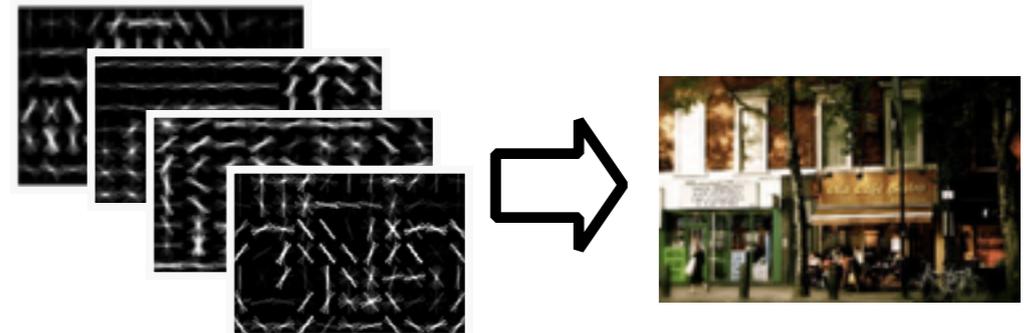
Immediate, scalable
object category detection

Motivation: Object Detection

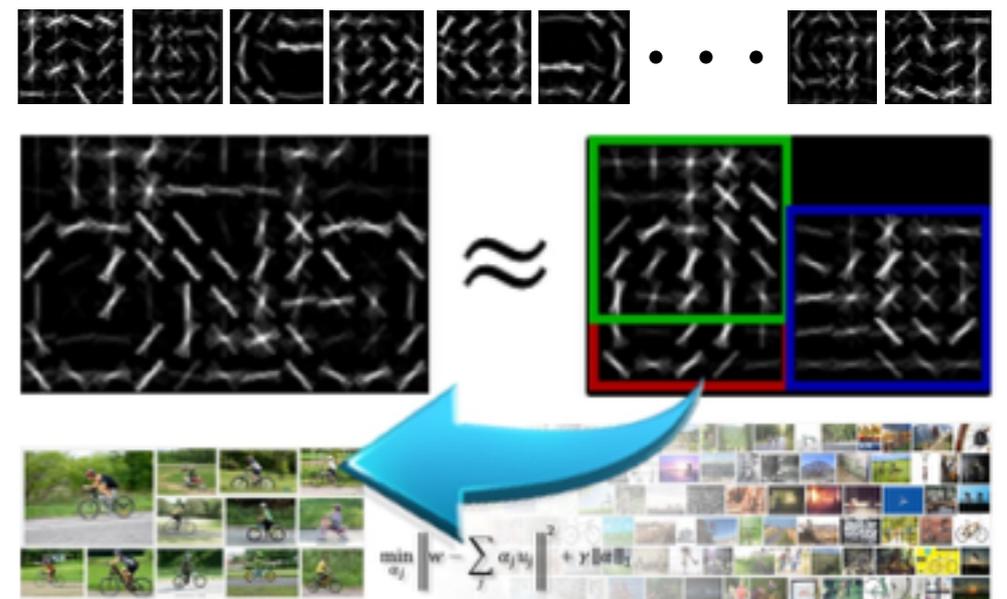
- Running *a detector* fast on a *single image* (Cascades, PQ, etc.) [Felzenszwalb-CVPR10, Vedaldi-CVPR12, Sadeghi-NIPS13].



- Running *multiple detectors* fast on a *single image* (Sparselets, etc.) [Song-ECCV12, Dean-CVPR13].



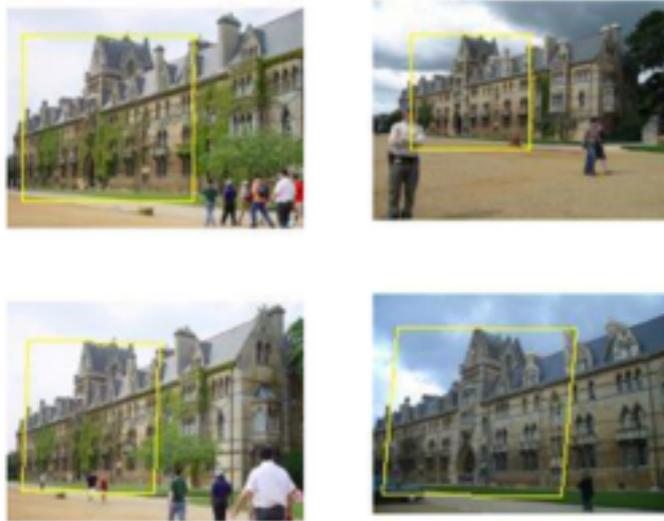
- Running *a detector* fast (*~1sec*) on a *large-scale image dataset*, similar to *Video Google* [Sivic03] but for *category detection*.



Large Scale Object Instance Retrieval



query

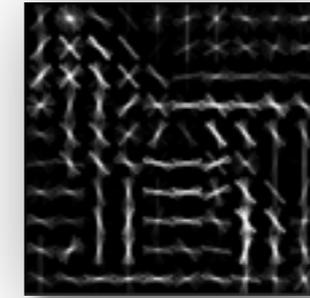


retrieval results

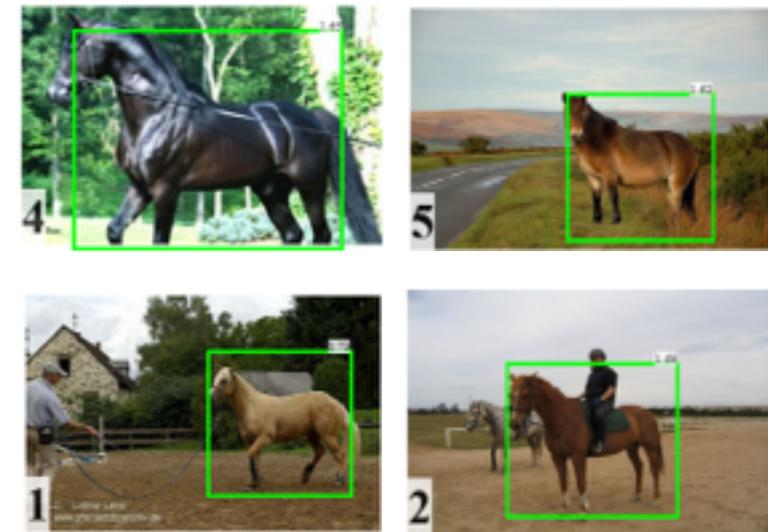
- Retrieve *instantly* ($< 1\text{sec}$)
- Robust to: *scale, viewpoint, lighting, partial occlusion*

versus

Large Scale Object Category Detection



query

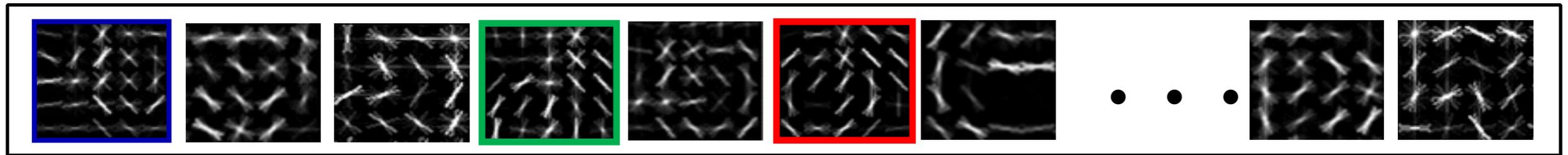


- Retrieve *instantly* ($\sim 1\text{sec}$)
- Robust to: *scale, viewpoint, lighting, partial occlusion and Intra-class variance*

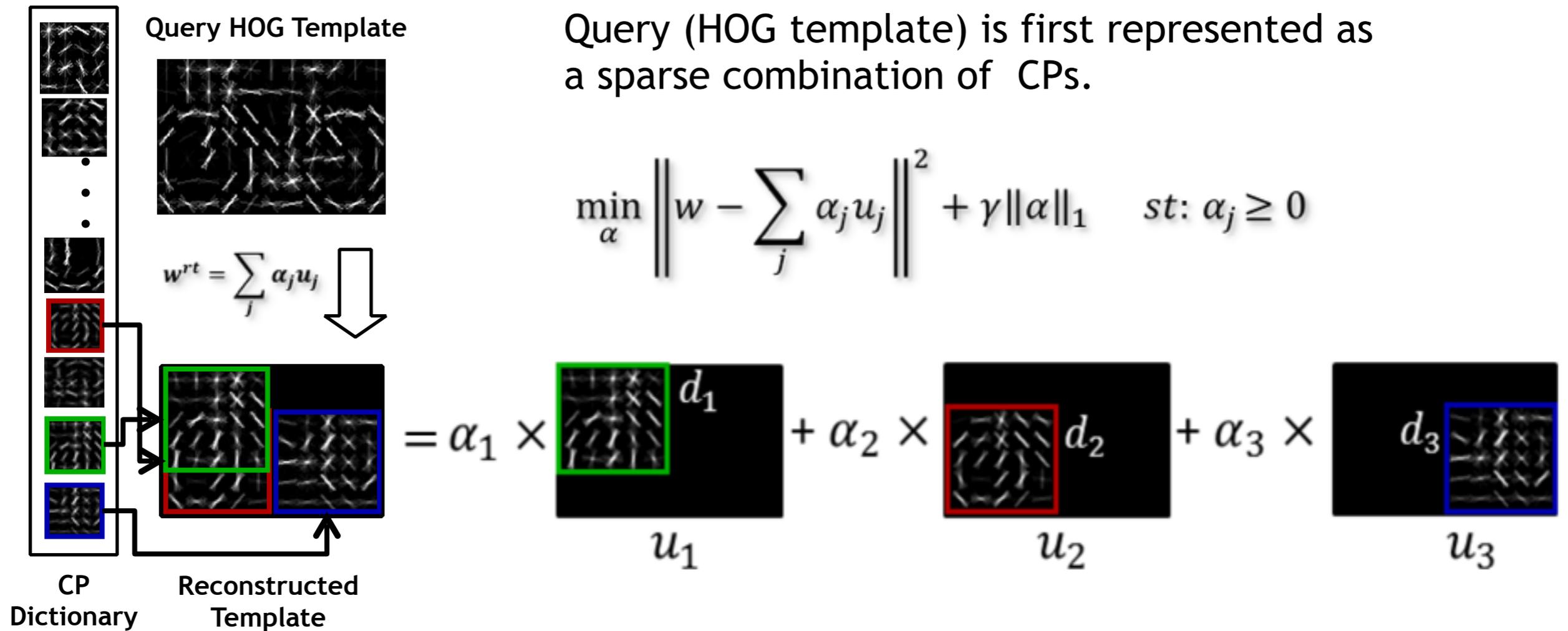
Overview

Uses the three stages of **Video Google** revamped for object category detection

Classifier Part (CP) Dictionary



Query Representation



Query (HOG template) is first represented as a sparse combination of CPs.

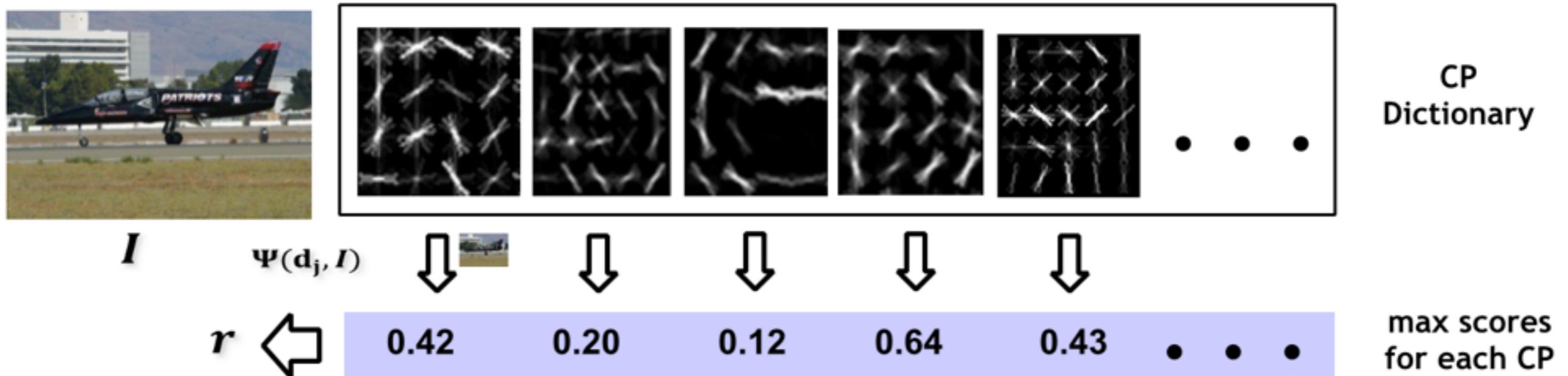
$$\min_{\alpha} \left\| w - \sum_j \alpha_j u_j \right\|^2 + \gamma \|\alpha\|_1 \quad st: \alpha_j \geq 0$$

α and the spatial layouts of CPs define the reconstructed template

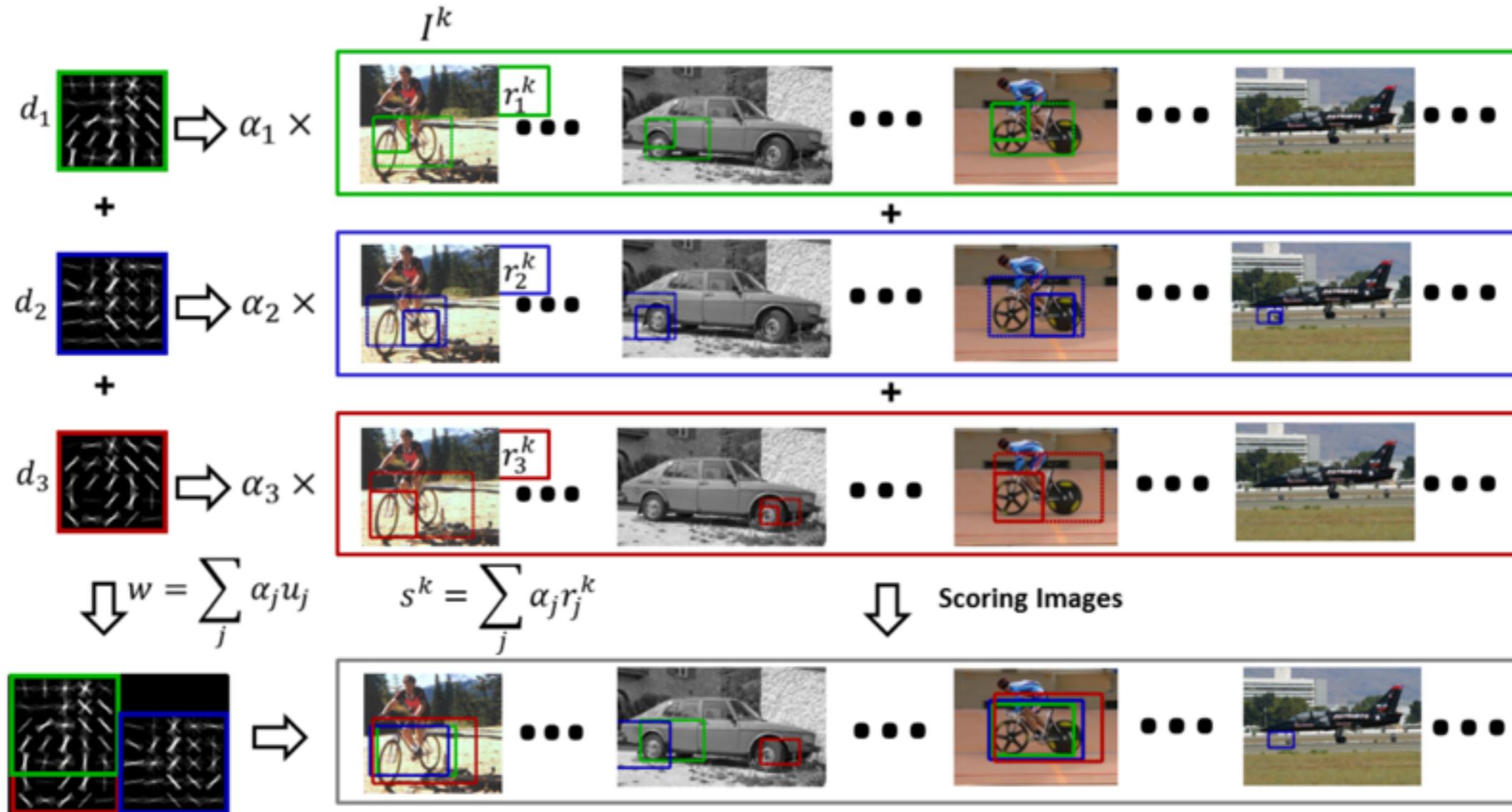
w^{rt}

Image Representation

Representation: Image I is represented with vector $\mathbf{r} = [\Psi(\mathbf{d}_1, I) \ \Psi(\mathbf{d}_2, I), \dots, \Psi(\mathbf{d}_M, I)]$, where the i^{th} component is the maximum response of the CP \mathbf{d}_i sliding over I .



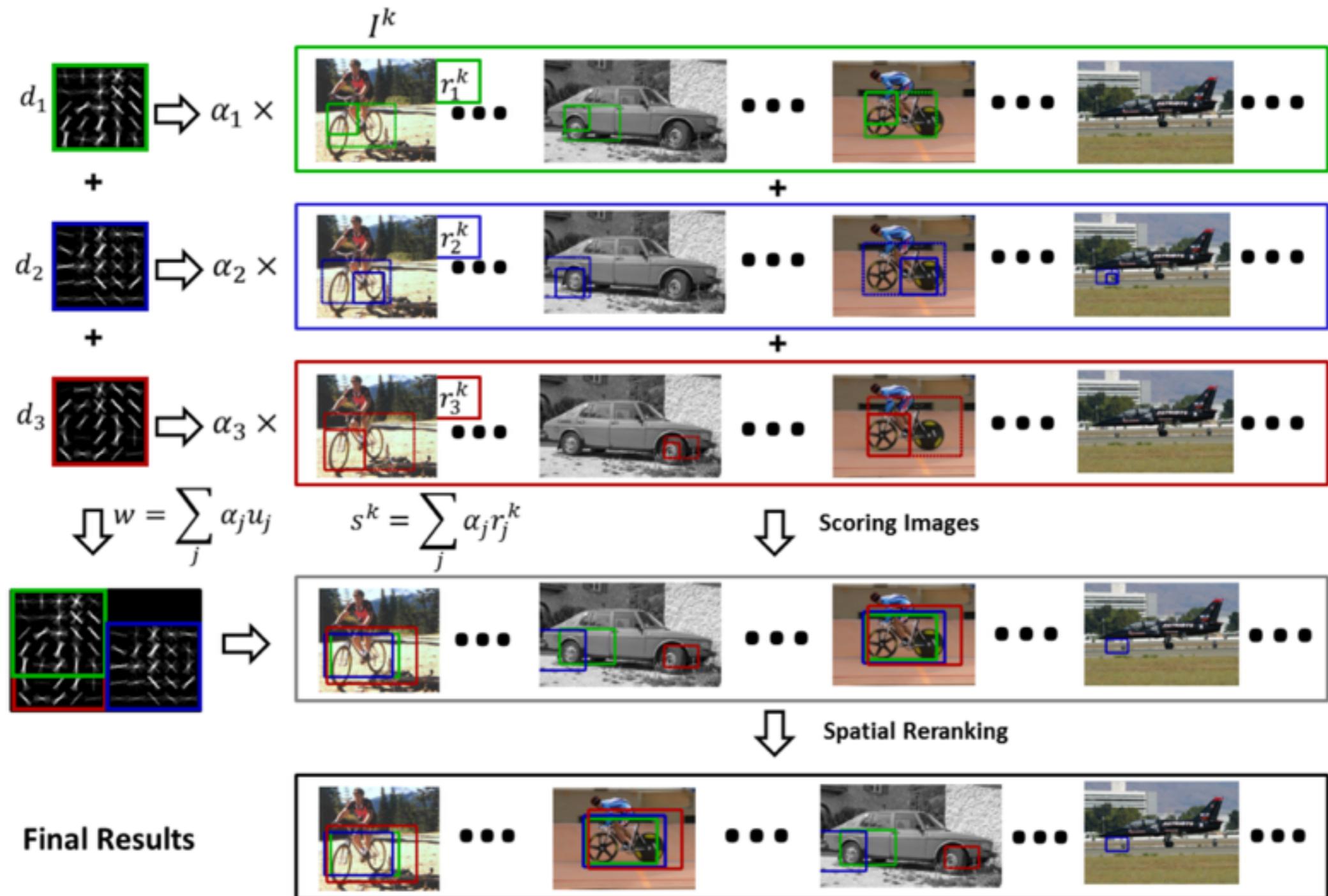
Shortlisting



Shortlisting: For the given template w^{rt} , images are shortlisted using the score $\alpha^T r$ which is an upper bound on the maximum score of w^{rt} obtained by sliding over I .

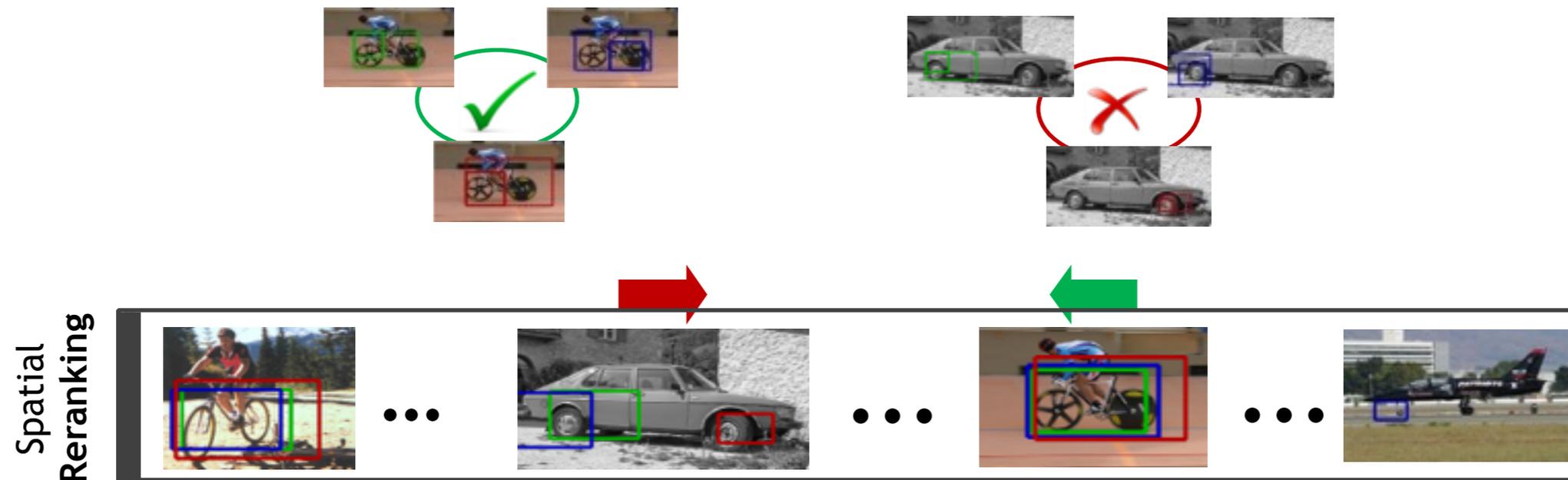
$$\alpha^T r = \sum_j \alpha_j \Psi(d_j, I) \geq \Psi(w^{rt}, I)$$

Spatial Reranking

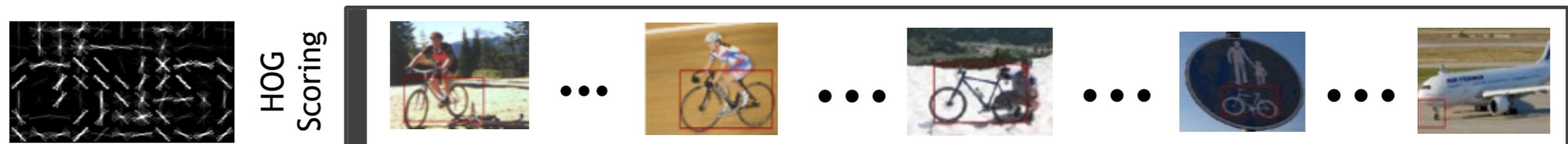


Reranking (Spatial + Original Template)

Shortlisted images are reranked via fast Hough-like voting of bounding box candidates suggested by each CP.

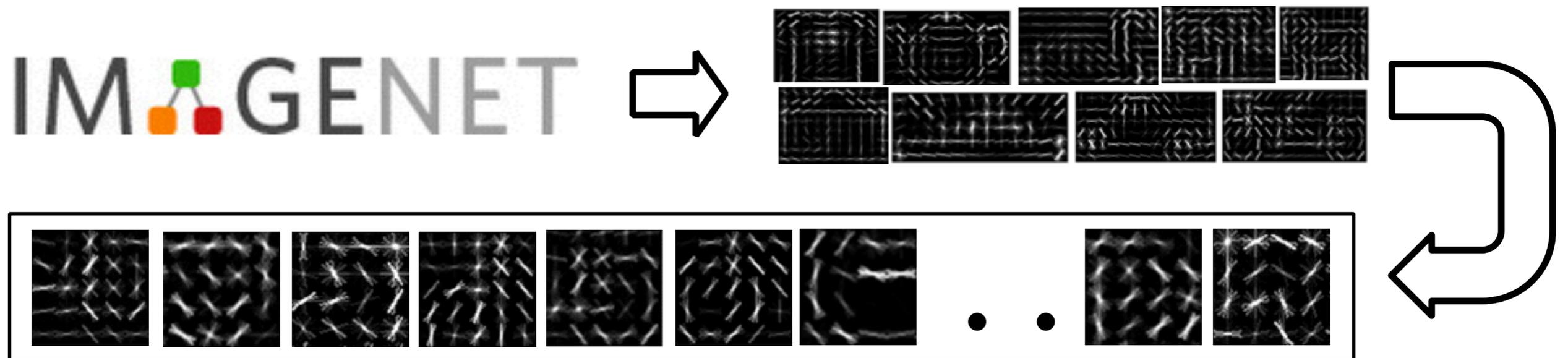


Retrieved bounding box candidates are re-scored using the original HOG template with fast and memory efficient PQ compression.



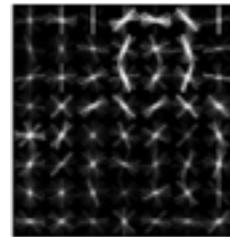
Dictionary & Dataset

10K dictionaries of sizes 3x3 - 7x7 HOG cells are extracted from DPMs trained from 1000 ImageNet categories.

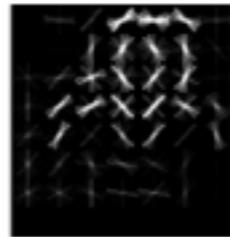


Tests are performed on PASCAL VOC07 test set (5K images) and validation sets (100K images) of ImageNet 2011 and 2012 challenges.

Detection Results



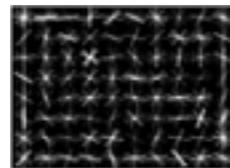
Person
Template



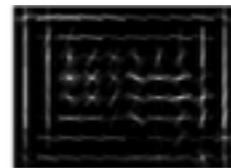
Reconstructed
Person Template



Top 3 retrievals



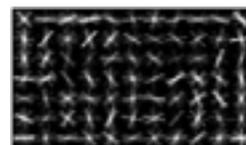
TV/Monitor
Template



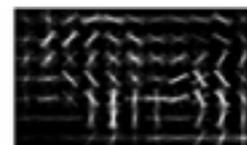
Reconstructed
TV/Monitor Template



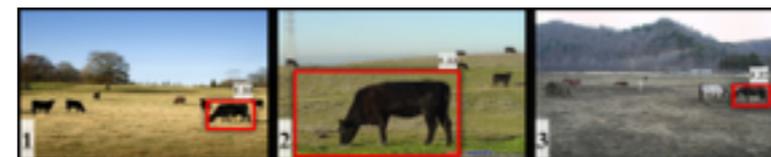
Top 3 retrievals



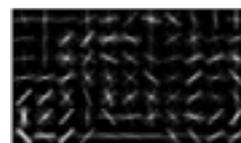
Cow
Template



Reconstructed
Cow Template



Top 3 retrievals



Motorbike
Template

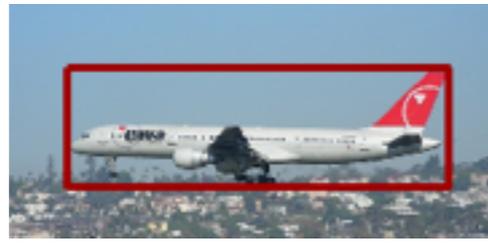


Reconstructed
Motorbike Template



Top 3 retrievals

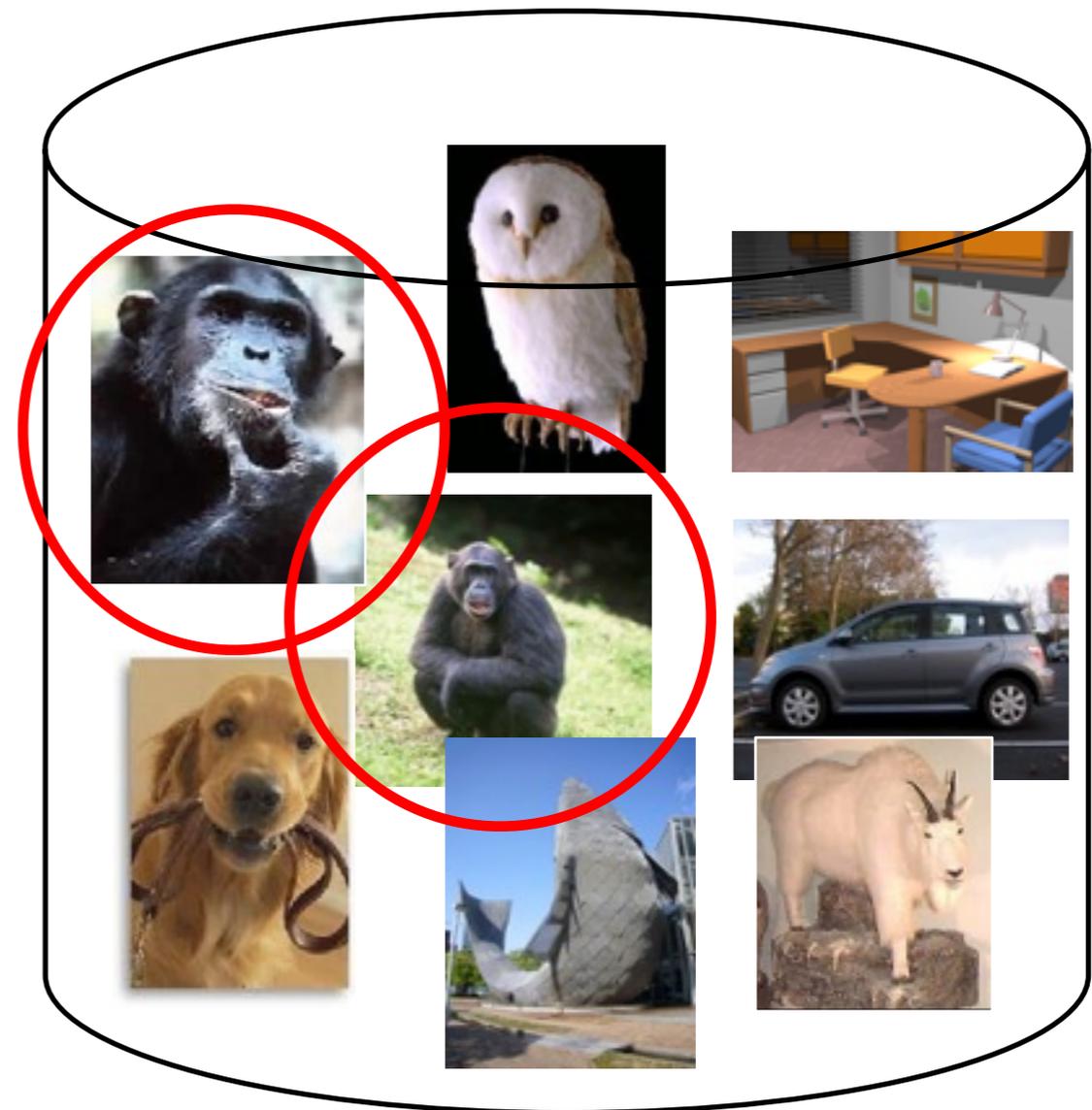
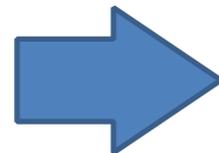
Exemplar SVM Results



Large Scale Image Search

Large Scale Image Search

- Find similar images in a large database



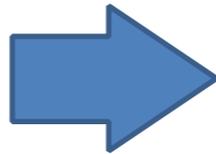
Fast & Accurate

Large Scale Image Search

Internet contains billions of images



Search the internet



The Challenge:

Need way of measuring similarity between images
(distance metric learning)

Needs to scale to Internet (How?)

Requirements for image search

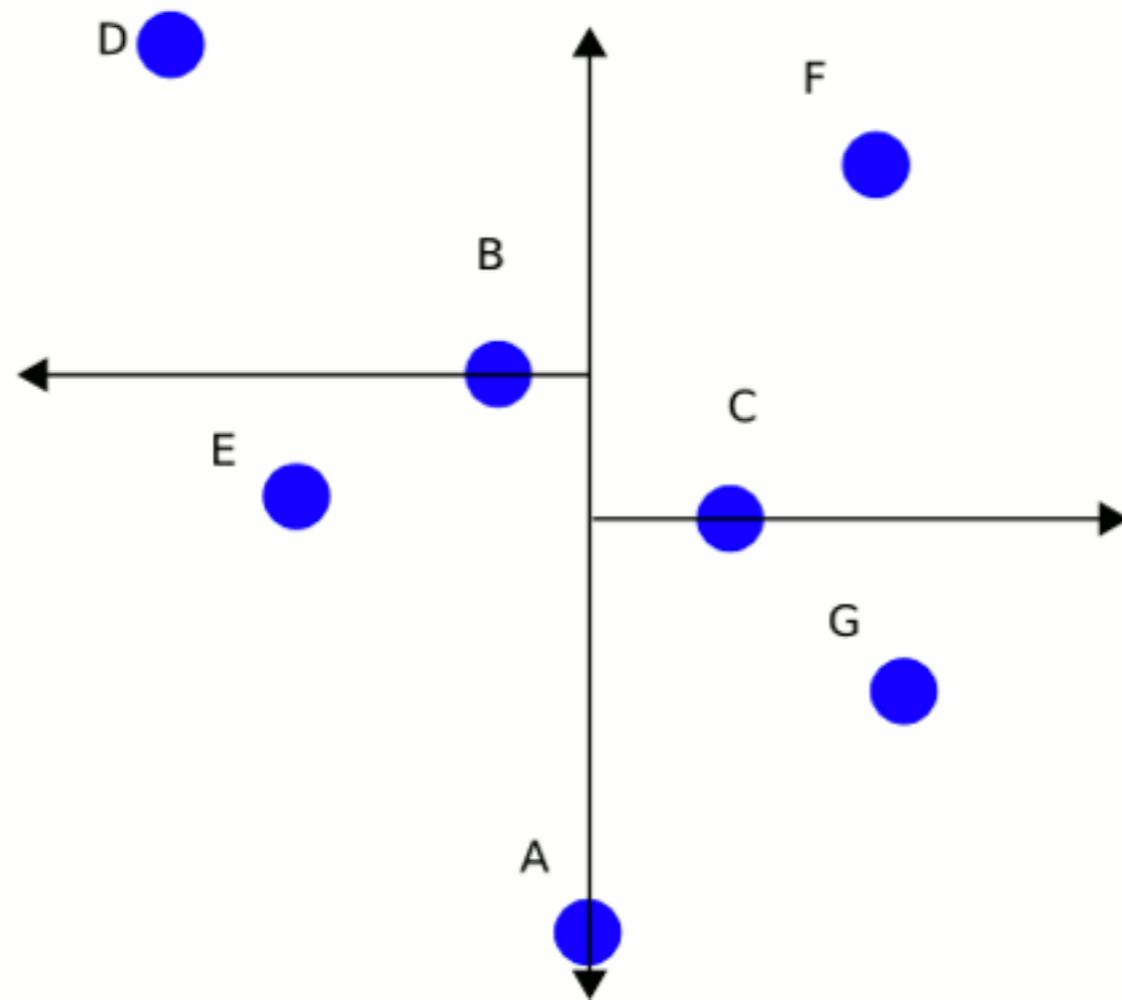
- Search must be both **fast, accurate** and **scalable to large data set**
- Fast
 - Kd-trees: tree data structure to improve search speed
 - Locality Sensitive Hashing: hash tables to improve search speed
 - Small code: binary small code (010101101)
- Scalable
 - Require very little memory, enabling their use on standard hardware or even on handheld devices
- Accurate
 - Learned distance metric

Categorization of existing large scale image search algorithms

- Tree Based Structure
 - Spatial partitions (i.e. **kd-tree**) and recursive hyper plane decomposition provide an efficient means to search low-dimensional vector data exactly.
- Hashing
 - **Locality-sensitive hashing** offers **sub-linear** time search by hashing highly similar examples together.
- Binary Small Code
 - **Compact binary code**, with a few hundred bits per image

Tree Based Structure

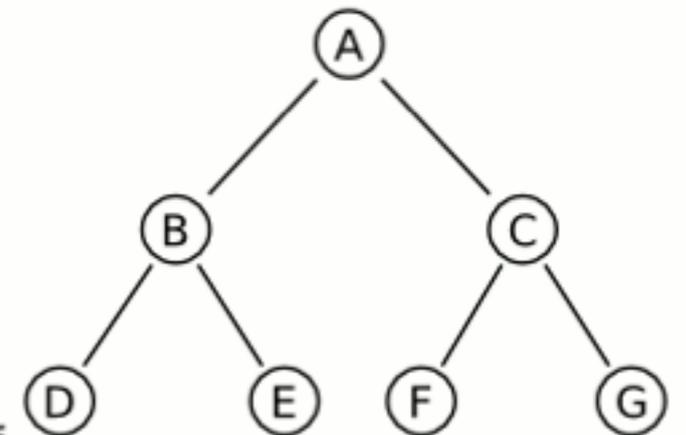
- Kd-tree
 - The kd-tree is a binary tree in which every node is a k-dimensional point



X-Splitting planes

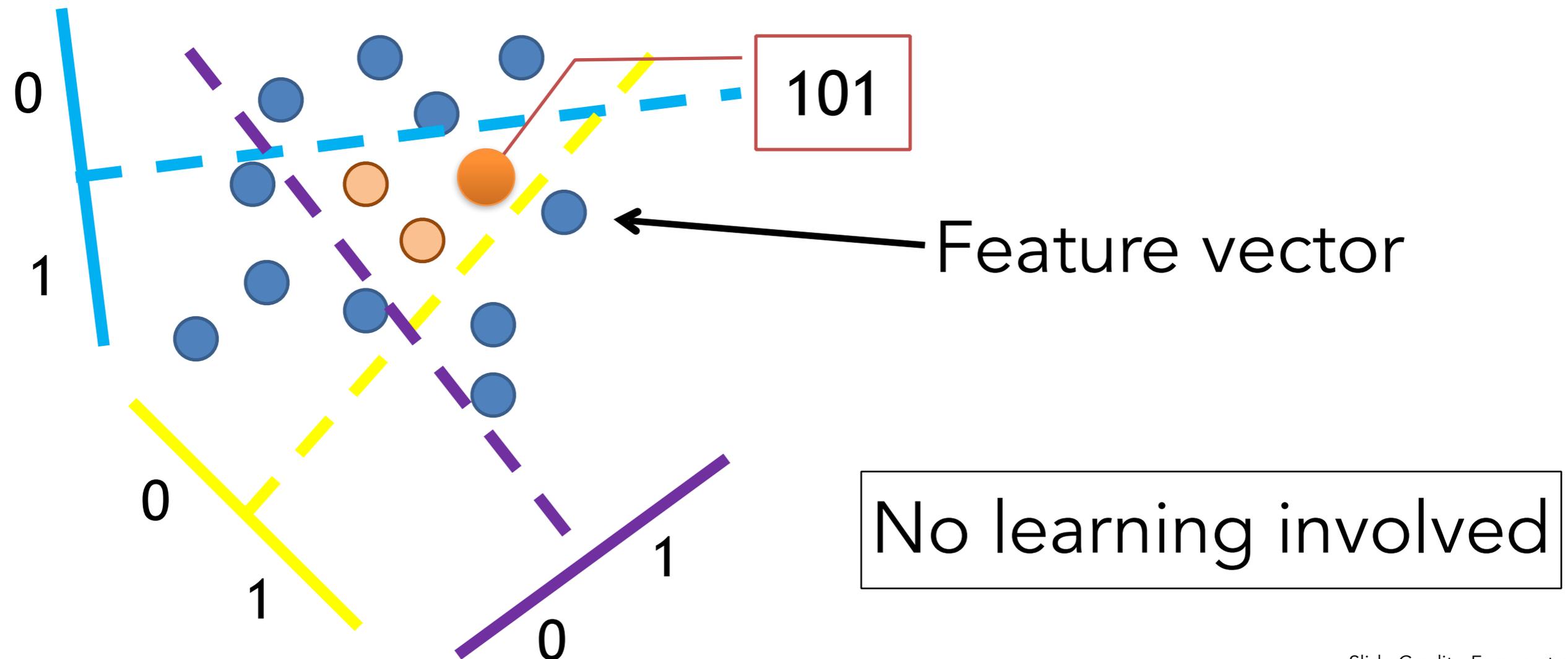
Y-Splitting planes

X-Splitting planes
not needed for leaf

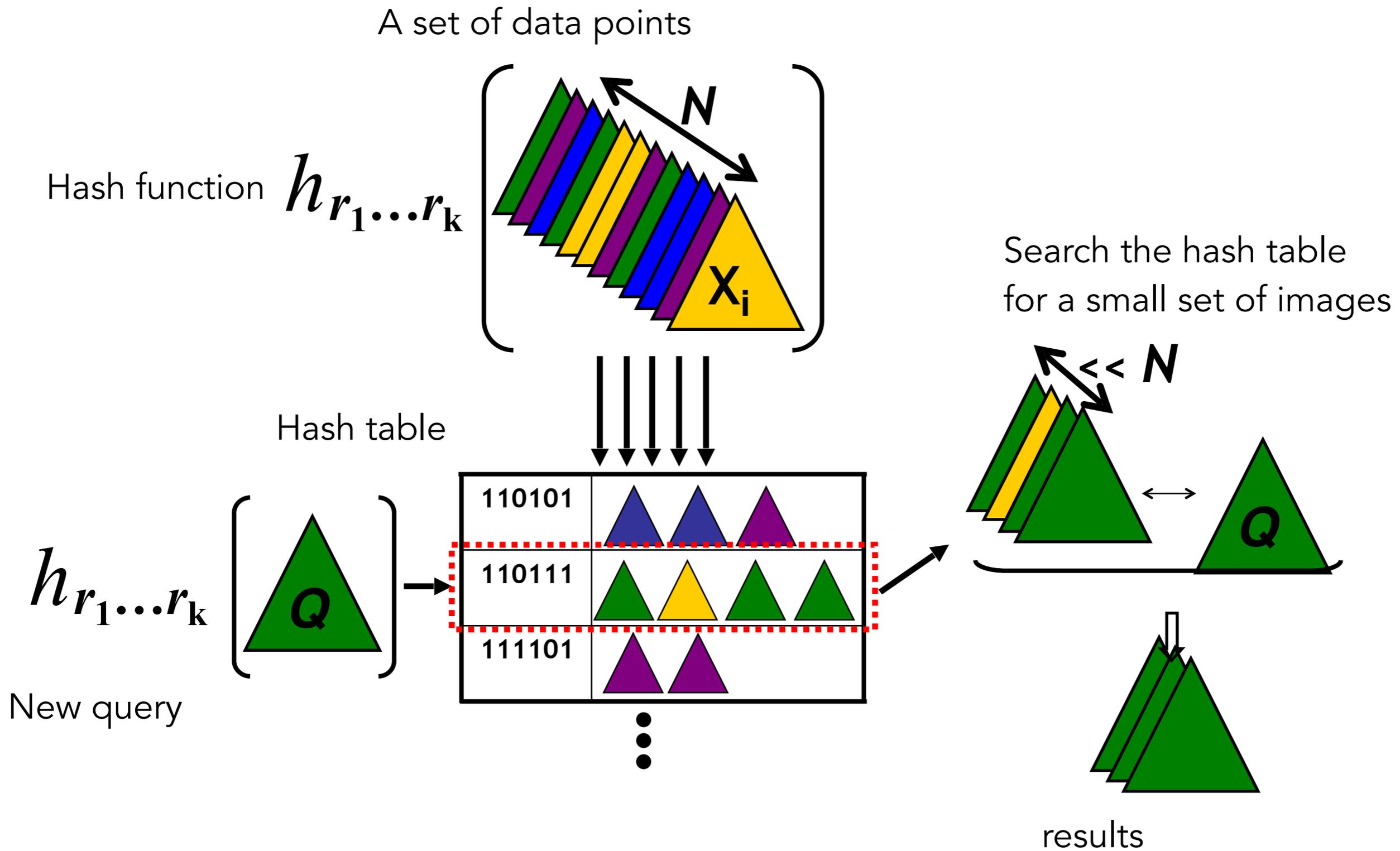


Locality Sensitive Hashing

- Take random projections of data $r^T x$
- Quantize each projection with few bits



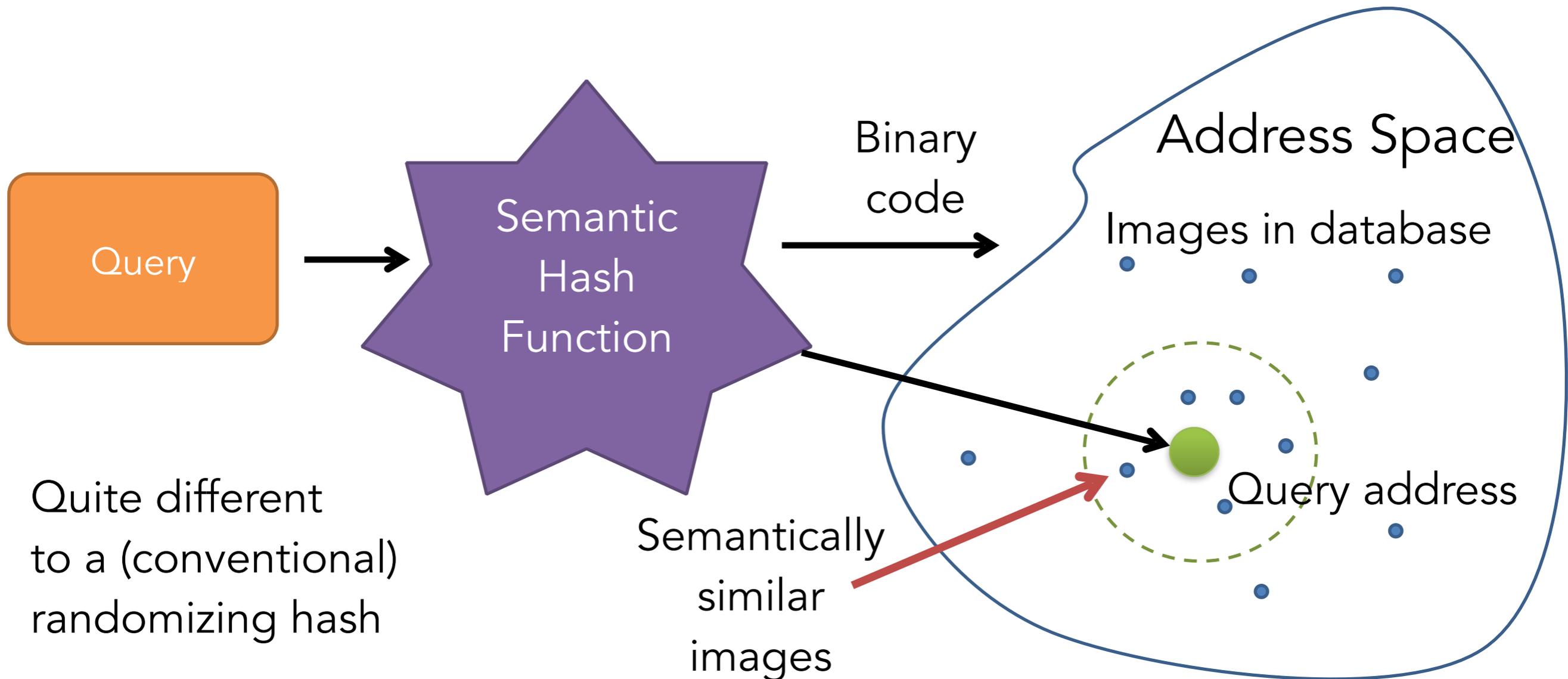
How to search from hash table?



Binary codes for images

- Want images with similar content to have similar binary codes
- Use Hamming distance between codes
 - Number of bit flips
 - E.g.: $\text{Ham_Dist}(10001010, 10001\mathbf{1}10) = 1$
 - $\text{Ham_Dist}(10001010, 1\mathbf{1}101\mathbf{1}10) = 3$
- Semantic Hashing [Salakhutdinov & Hinton, 2007]
 - Text documents

Semantic Hashing



- Find neighbors by exploring Hamming ball around query address
- Lookup time depends on radius of ball, NOT on # data points

Compact Binary Codes

- Google has few billion images (10^9)
- PC has ~ 10 Gbytes (10^{11} bits)
- Codes must fit in memory (disk too slow)

Budget of **10^2 bits/image**

- 1 Megapixel image is 10^7 bits
- 32x32 color image is 10^4 bits

Semantic hash function must also reduce dimensionality

RBM architecture

- Network of binary stochastic units
- Hinton & Salakhutdinov, Science 2006

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i \in \text{visible}} b_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij}$$

Parameters: Weights w Biases b

Convenient conditional distributions:

$$p(h_j = 1 | \mathbf{v}) = \sigma(b_j + \sum_i w_{ij} v_i)$$

$$p(v_i = 1 | \mathbf{h}) = \sigma(b_i + \sum_j w_{ij} h_j)$$

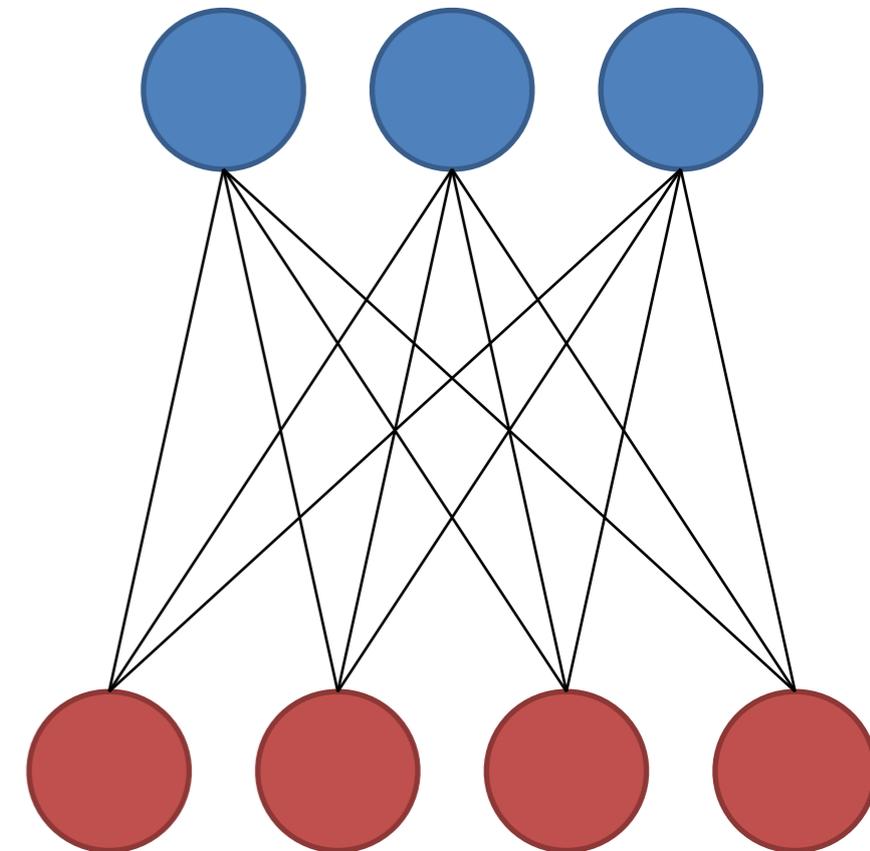
$\sigma(x) = 1/(1 + e^{-x})$, the logistic function

Learn weights and biases using
Contrastive Divergence

Hidden units: h

Symmetric
weights w

Visible units: v



Examples of LabelMe retrieval

12 closest neighbors under different distance metrics

