

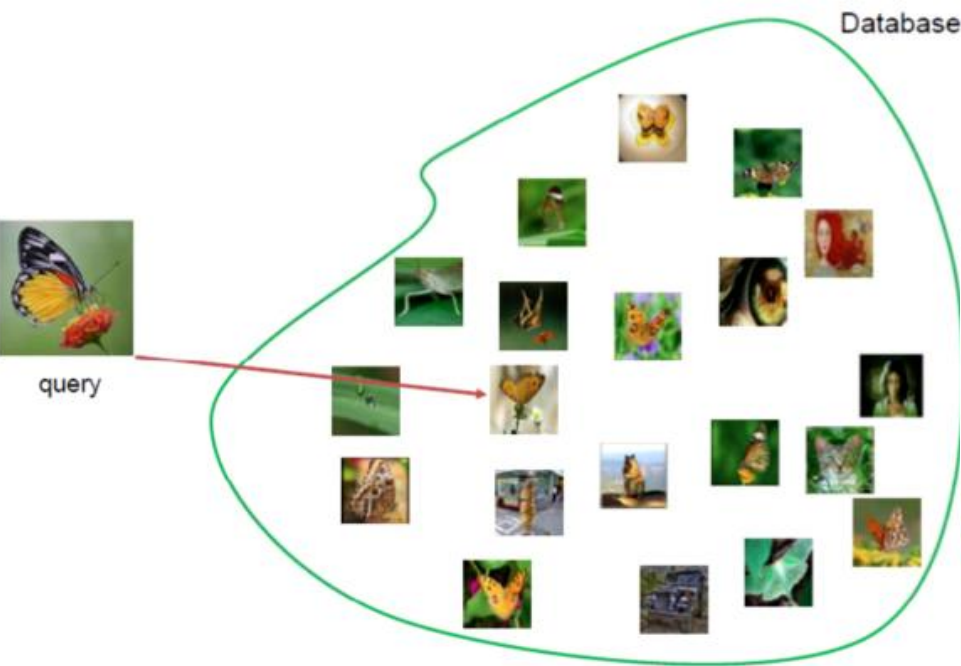
# Supervised Hashing for Image Retrieval via Image Representation Learning

Rongkai Xia, Yan Pan, Cong Liu (**Sun Yat-Sen University**)

Hanjiang Lai, Shuicheng Yan (**National University of Singapore**)

# Finding Similar Images

**Task:** given a query image, find its nearest neighbors in an image database.



search



# Challenges

Query image



features

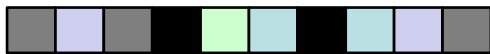
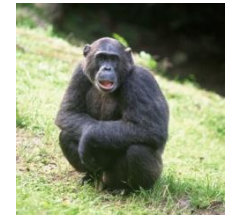
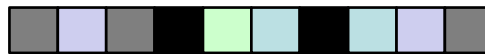


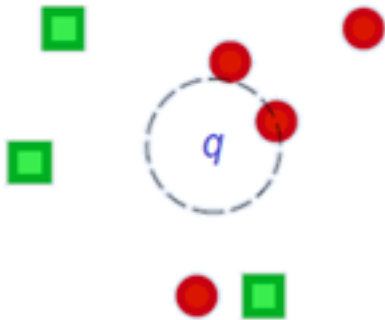
Image database



features



features



**Challenge:** how to efficiently search over millions or billions of data?

- e.g., if each image is represented by a **512-dim** GIST vector, one needs **20G** memory to store **10 million** images.

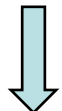
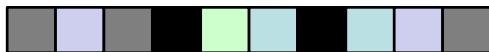
Nearest neighbors search

# Similarity-Preserving Hashing

Query image



features

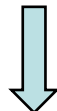
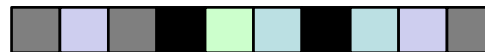


1001010

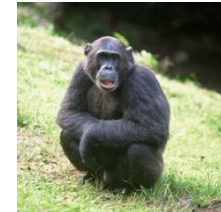
Image database



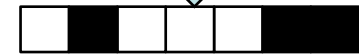
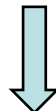
features



1001000



feature



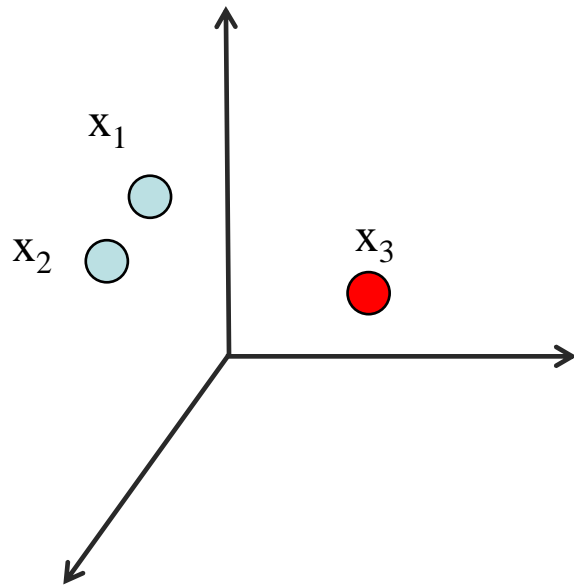
0100011

Images are represented by **binary codes**.

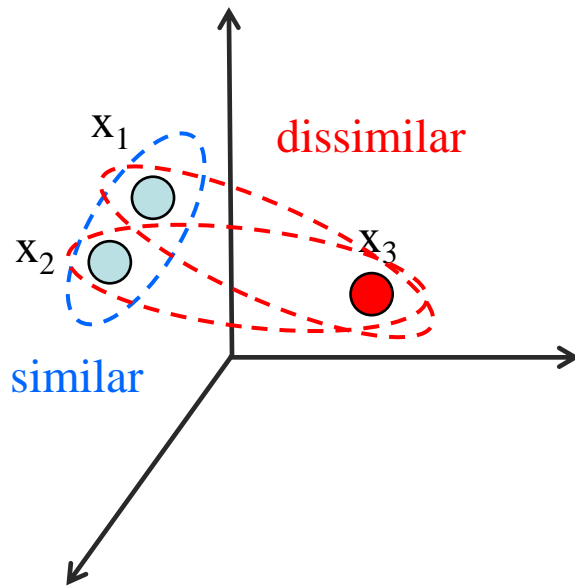
- Efficient retrieval via bitwise operations.
- Space-efficient storage
  - e.g., 10 million images in 32M memory, providing each image in 128-bit code

Key question: how to preserve similarities?

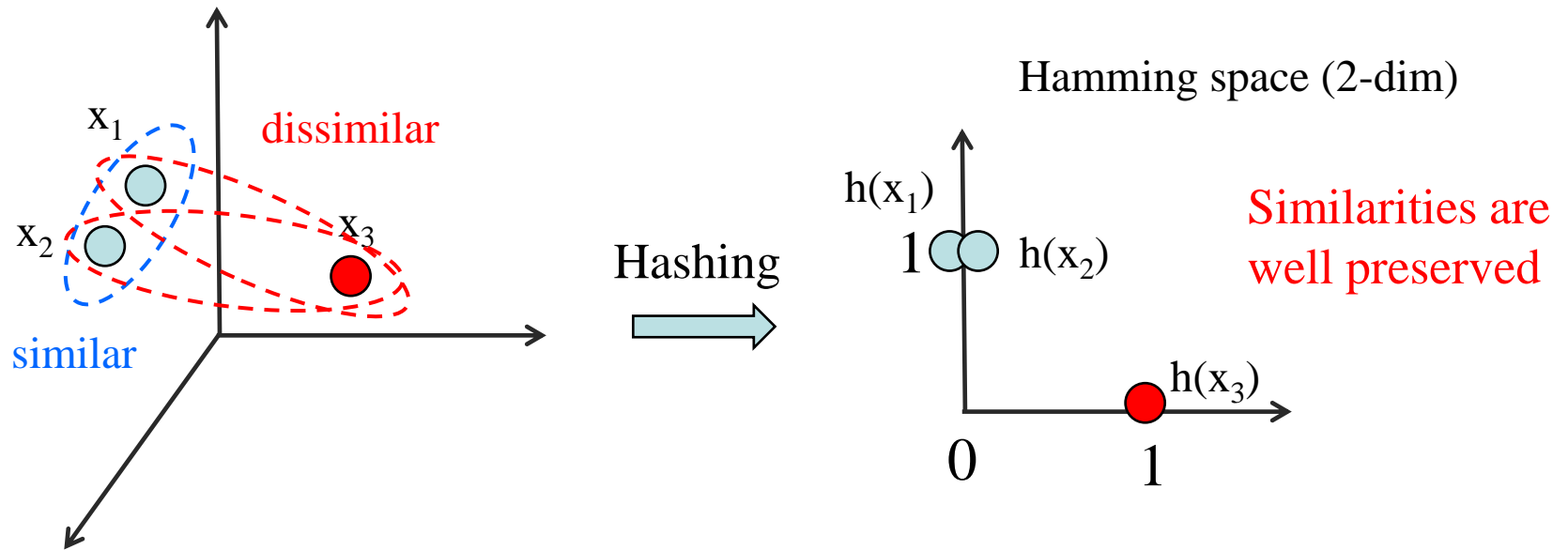
# Similarity-Preserving Hashing



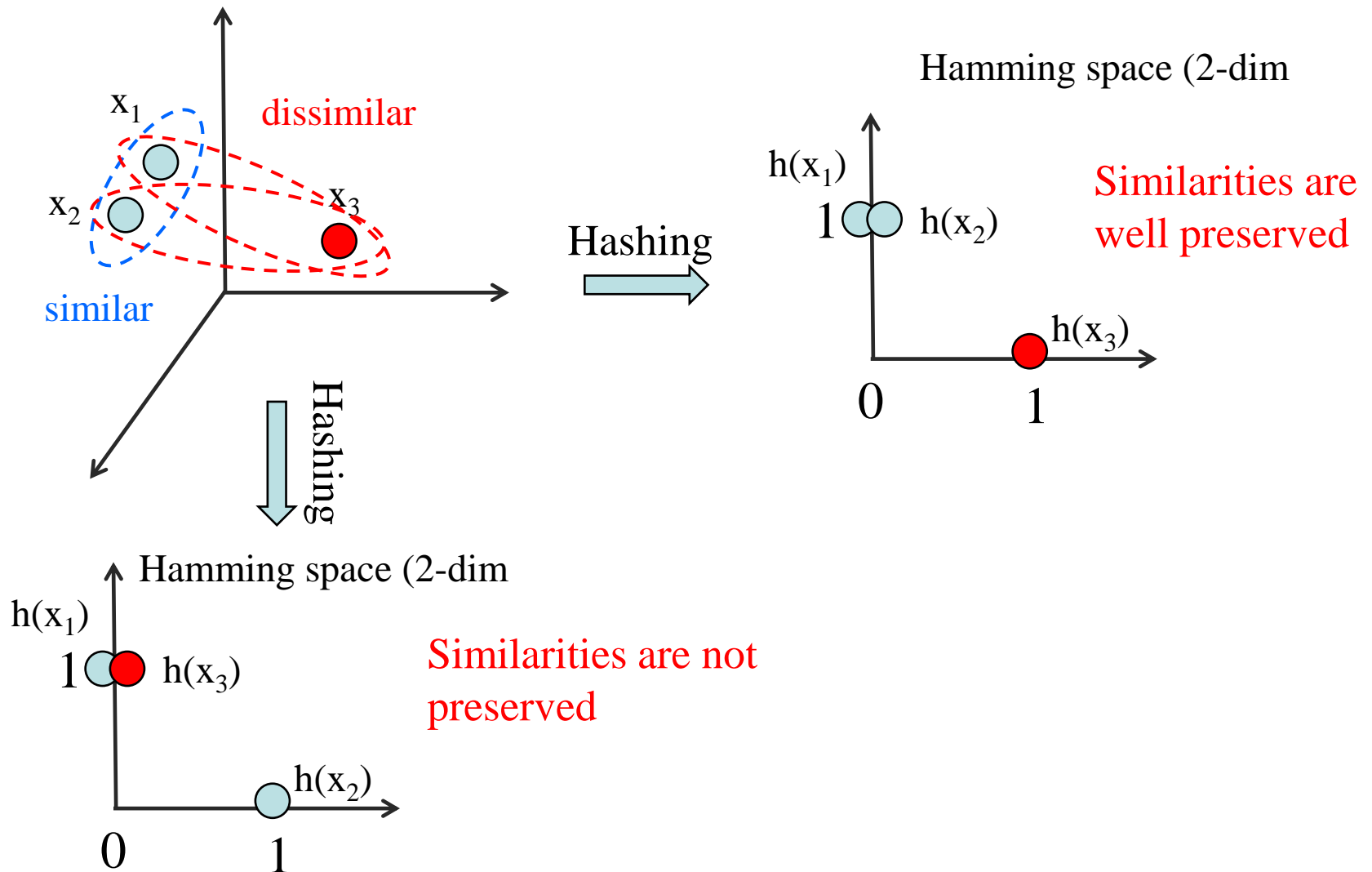
# Similarity-Preserving Hashing



# Similarity-Preserving Hashing



# Similarity-Preserving Hashing





# Related Work

Long codes are needed to preserve similarities.

Unsupervised Hashing

**LSH**[Gionis et al. VLDB,1999]  
**KLSH**[Kulis and Grauman. PAMI,2012]  
**SH**[Weiss and Torralba. NIPS,2008]  
**ITQ**[Gong and Lazebnik. CVPR,2011]

Semi-supervised Hashing

**SSH**[Wang et al. CVPR,2010]

Supervised Hashing

**MLH**[Norouzi and Blei. ICML,2011]  
**BRE**[Kulis and Darrell. NIPS,2009]  
**KSH**[Liu et al. CVPR,2012]  
**TSH**[Lin et al. ICCV,2013]

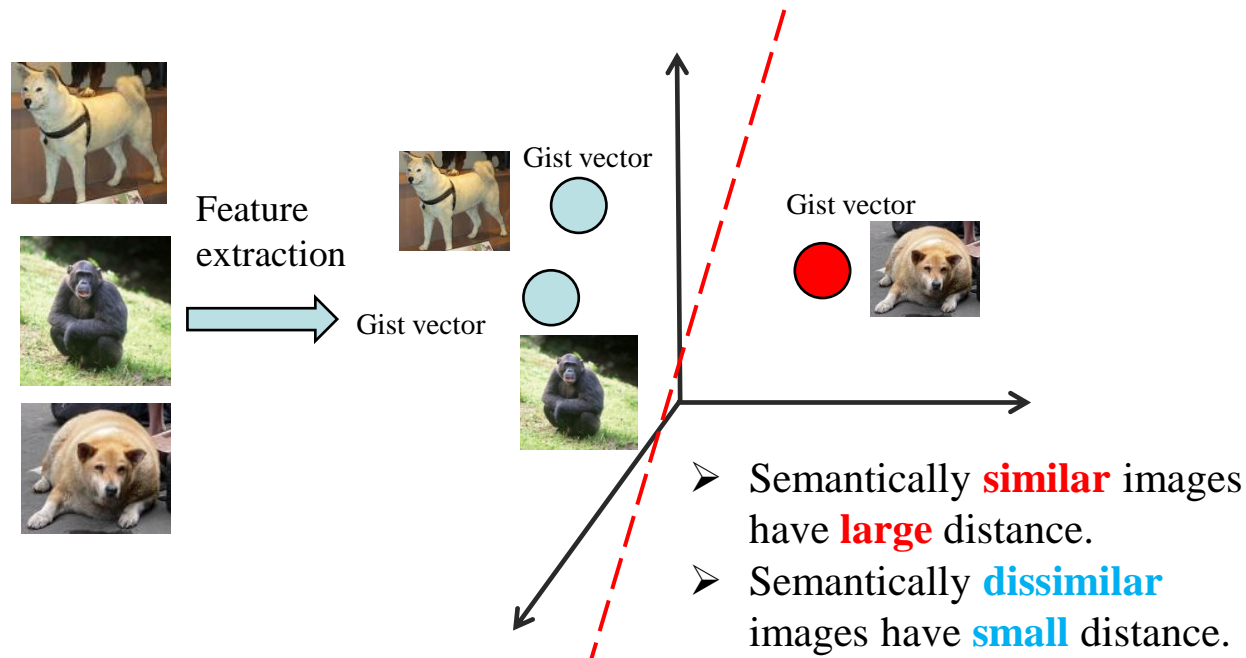
Learn compact codes by using label information.

# Motivation

- In most existing methods, each image is firstly encoded by a vector of some hand-crafted visual descriptor (e.g., GIST, BoW, SIFT)
- **Concern:** the chosen hand-crafted visual features do not necessarily guarantee to accurately preserve the semantic similarities of image pairs.
  - e.g., a pair of semantically similar/dissimilar images may not have feature vectors with relatively small/large Euclidean distance.

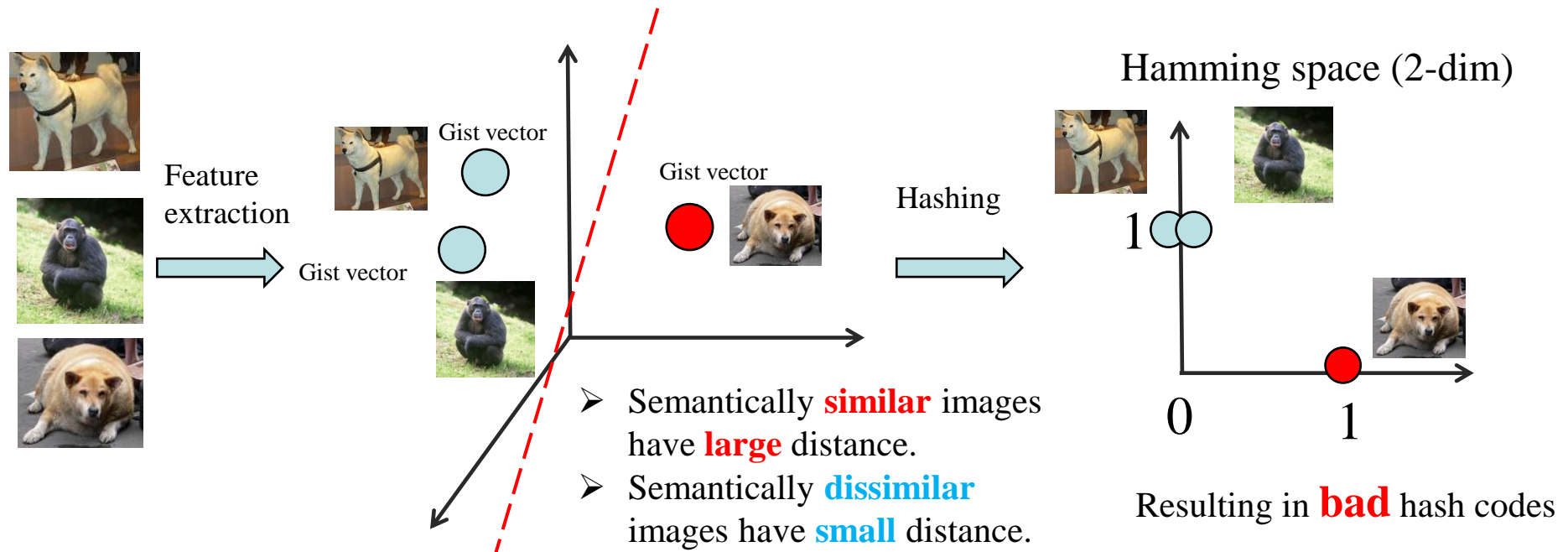
# Motivation

- In most existing methods, each image is firstly encoded by a vector of some hand-crafted visual descriptor (e.g., GIST, BoW, SIFT) .
- **Concern:** the chosen hand-crafted visual features do not necessarily guarantee to accurately preserve the semantic similarities of image pairs.
  - e.g., a pair of semantically similar/dissimilar images may not have feature vectors with relatively small/large Euclidean distance.



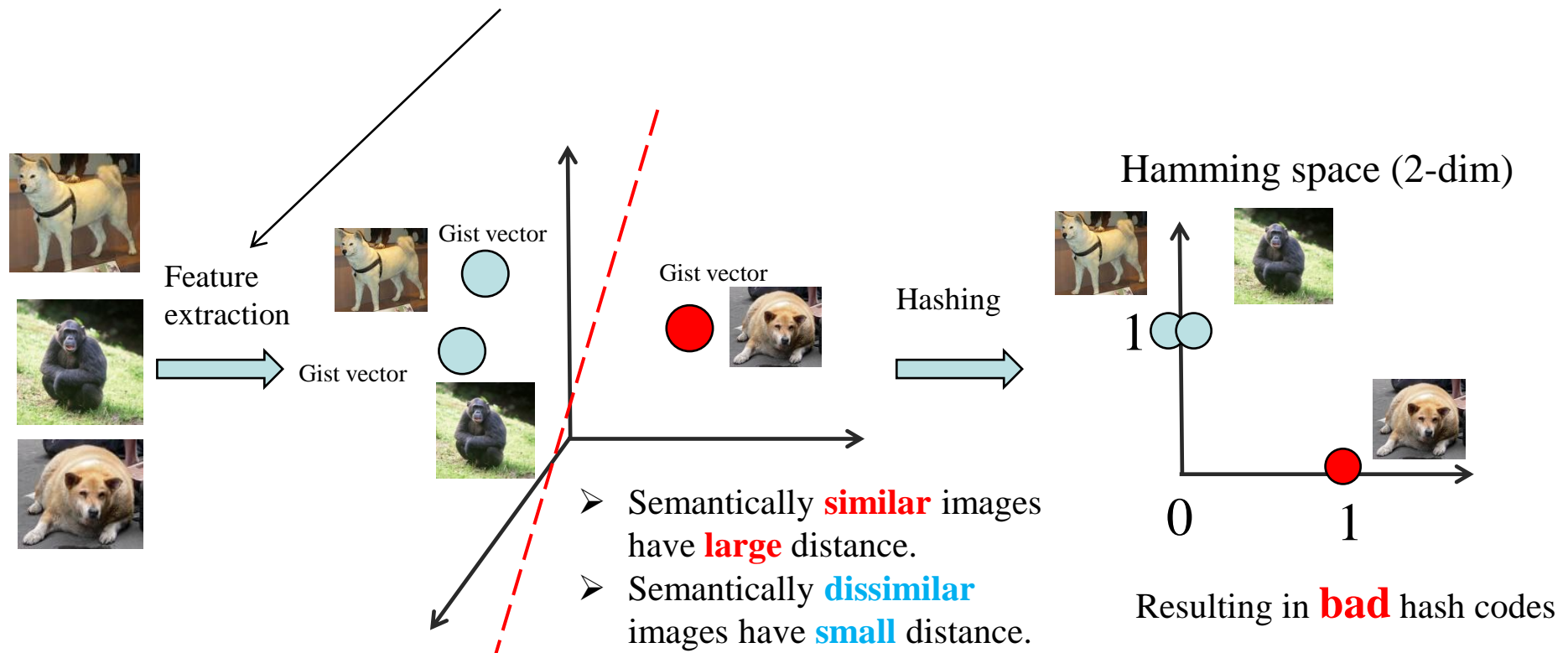
# Motivation

- In most existing methods, each image is firstly encoded by a vector of some hand-crafted visual descriptor (e.g., GIST, BoW, SIFT)
- **Concern:** the chosen hand-crafted visual features do not necessarily guarantee to accurately preserve the semantic similarities of image pairs.
  - e.g., a pair of semantically similar/dissimilar images may not have feature vectors with relatively small/large Euclidean distance.



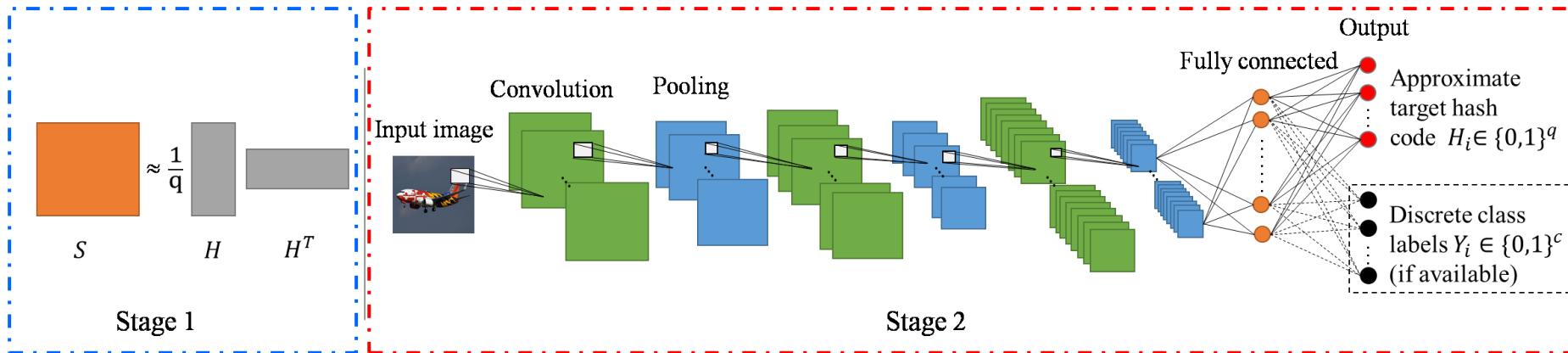
# Motivation

A useful image representation is **important** in hash learning process.



# The Proposed Approach

## Two-stage framework



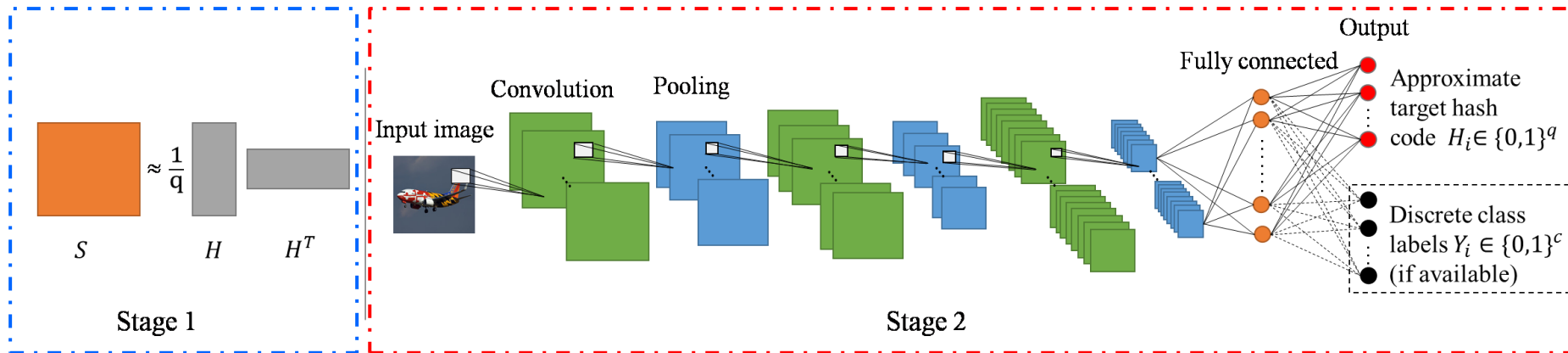
1. **Learn approximate hash codes for the training samples**, i.e., the pairwise similarity matrix  $S$  is decomposed into a product  $S = \frac{1}{q} H H^T$  where the  **$i$ th row in  $H$**  is the approximate hash code for the  **$i$ th training image**

$$\min_H \|S - \frac{1}{q} H H^T\|_F^2 \text{ s.t. } H \in [-1, 1]^{n \times q}. S_{ij} = \begin{cases} +1, & I_i, I_j \text{ are semantically similar} \\ -1, & I_i, I_j \text{ are semantically dissimilar.} \end{cases}$$

2. By using the learnt  **$H$**  and the raw image pixels as input, learn **image representation** and **hash functions** via *deep convolutional neural networks*.

# The Proposed Approach

## Two-stage framework



## Stage 1

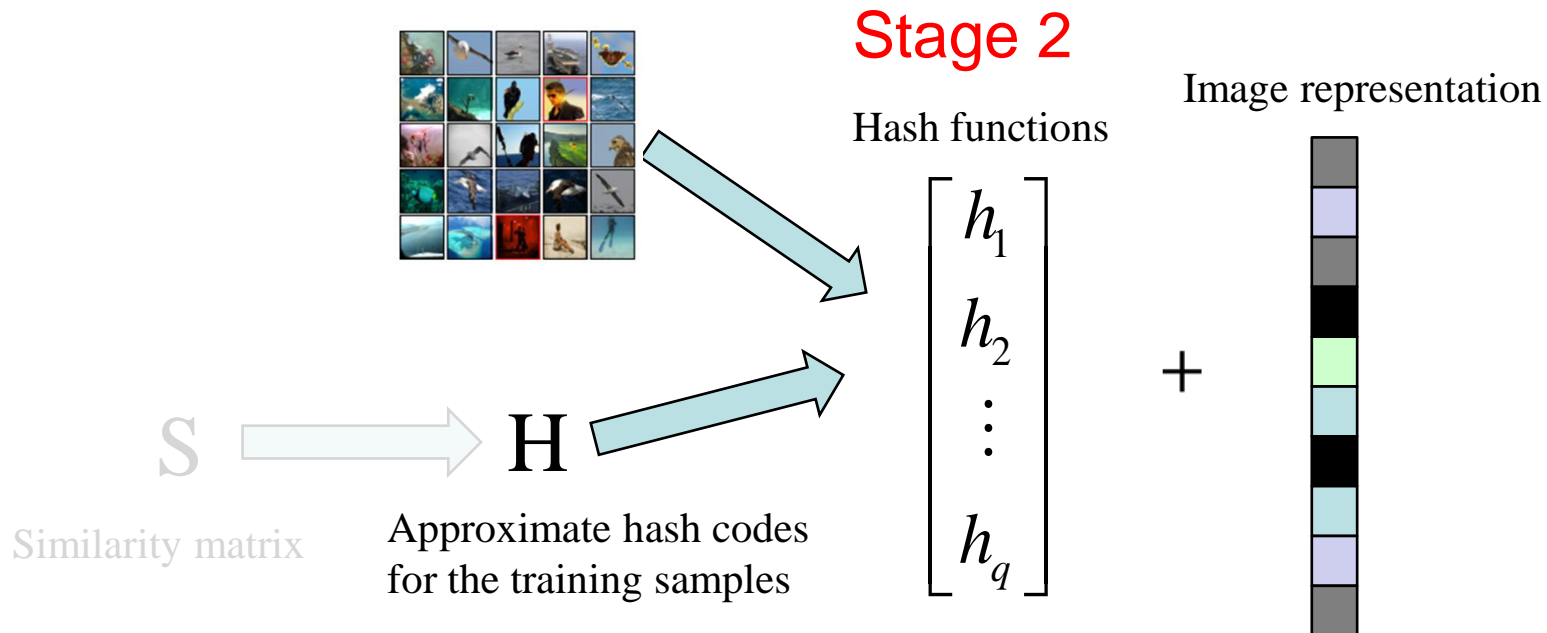
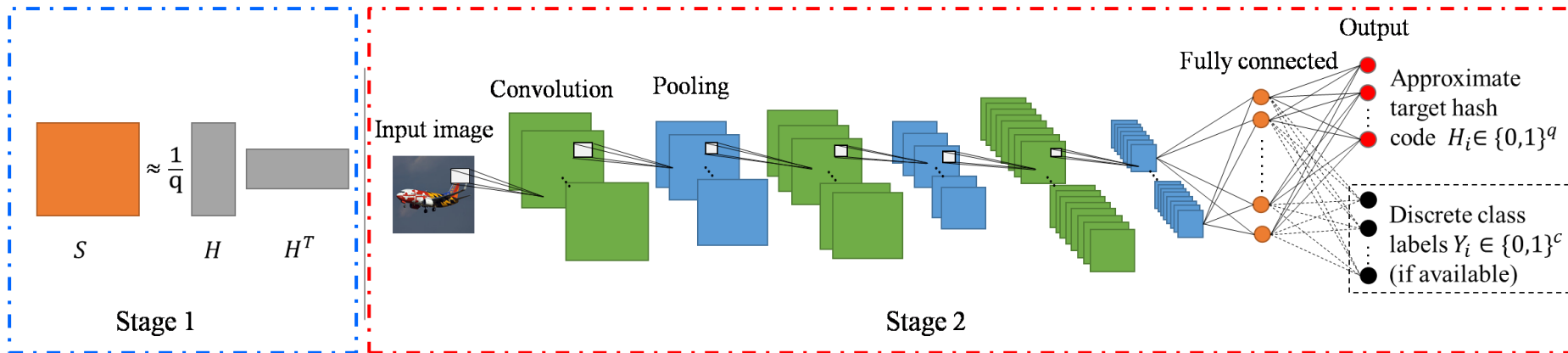


Similarity matrix

Approximate hash codes  
for the training samples

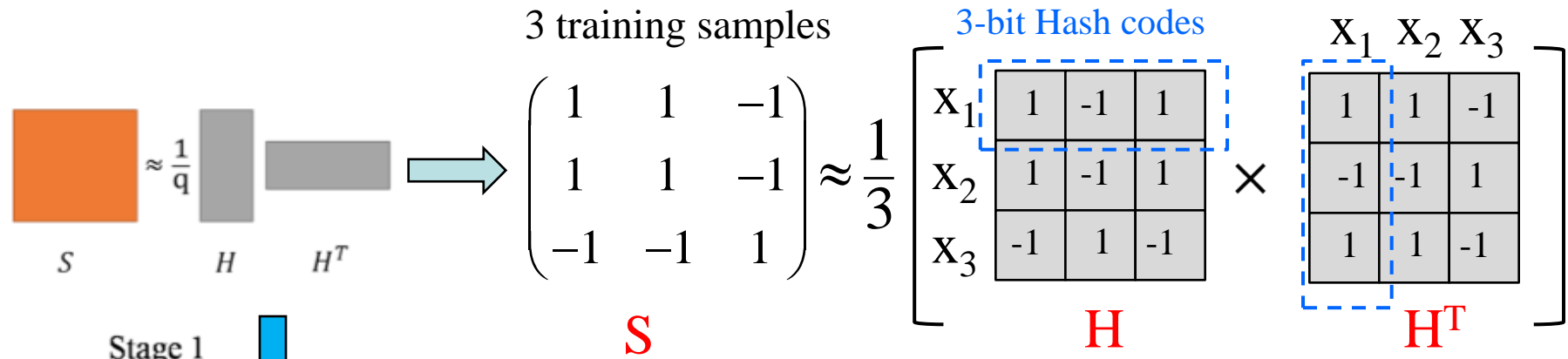
# The Proposed Approach

## Two-stage framework





# Stage 1: Learning Approximate Codes



$$\min_H \sum_{i=1}^n \sum_{j=1}^n \left( S_{ij} - \frac{1}{q} H_i \cdot H_j^T \right)^2$$

$$= \min_H \left\| S - \frac{1}{q} H H^T \right\|_F^2$$

$$\text{subject to : } H \in \{-1, 1\}^{n \times q}$$

The Hamming distance of two hash codes has one-to-one correspondence to the inner product of these two codes.

If images  $i$  and  $j$  are **similar** (**dissimilar**), the inner product of their approximate codes should be **large** (**small**).

$$\min_H \left\| S - \frac{1}{q} H H^T \right\|_F^2 \text{ s.t. } H \in [-1, 1]^{n \times q}$$

**relaxation**

# Optimization

$$\min_H \|S - \frac{1}{q} H H^T\|_F^2 \text{ s.t. } H \in [-1, 1]^{n \times q}.$$

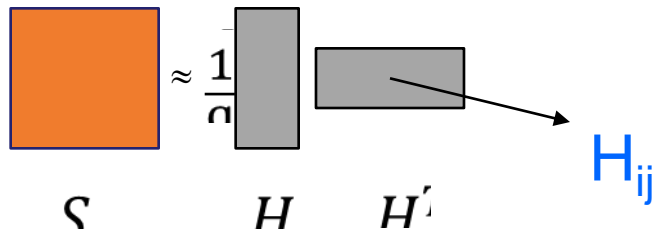
- **Algorithm:** random coordinate descent using Newton directions
  1. Randomly select an entry  $H_{ij}$  in  $H$  to update while keeping other entries fixed
  2. Approximate the objective function by second-order Taylor expansion w.r.t.  $H_{ij}$   
Calculate a step-size  $d$  and update  $H_{ij}$  by  $H_{ij} \leftarrow H_{ij} + d$
  3. Repeat 1 and 2 until stopping criterion is satisfied

# Optimization

$$\min_H \left\| S - \frac{1}{q} H H^T \right\|_F^2 \quad \text{s.t. } H \in [-1, 1]^{n \times q}.$$

- **Algorithm:** random coordinate descent using Newton directions

1. Randomly select an entry  $H_{ij}$  in  $H$  to update while keeping other entries fixed



2. Approximate the objective function by second-order Taylor expansion w.r.t.  $H_{ij}$   
Calculate a step-size  $d$  and update  $H_{ij}$  by  $H_{ij} \leftarrow H_{ij} + d$
3. Repeat 1 and 2 until stopping criterion is satisfied

# Optimization

$$\min_H \|S - \frac{1}{q} H H^T\|_F^2 \text{ s.t. } H \in [-1, 1]^{n \times q}.$$

- **Algorithm:** random coordinate descent using Newton directions

1. Randomly select an entry  $H_{ij}$  in  $H$  to update while keeping other entries fixed
2. Approximate the objective function by second-order Taylor expansion w.r.t.  $H_{ij}$   
Calculate a step-size  $d$  and update  $H_{ij}$  by  $H_{ij} \leftarrow H_{ij} + d$

$$\begin{aligned} \min_{H_{ij}} g(H_{ij}) &= \sum_{l=1}^n \sum_{k=1}^n (H_{lj} H_{kj} - R_{lk})^2 \\ &= (H_{ij}^2 - R_{ii})^2 + 2 \sum_{k \neq i} (H_{ij} H_{kj} - R_{ik})^2 + \text{constant} \end{aligned} \quad \Longrightarrow \quad \begin{aligned} \min_d g(H_{ij} + d) \\ \text{subject to : } -1 \leq H_{ij} + d \leq 1 \end{aligned}$$

subject to :  $H_{ij} \in [-1, 1]$

$$\Longrightarrow g(H_{ij} + d) \approx g(H_{ij}) + g'(H_{ij})d + \frac{1}{2} g''(H_{ij})d^2 \Longrightarrow d = \max(-1 - H_{ij}, \min(-\frac{g'(H_{ij})}{g''(H_{ij})}, 1 - H_{ij}))$$

3. Repeat 1 and 2 until stopping criterion is satisfied

# Optimization

$$\min_H \|S - \frac{1}{q} H H^T\|_F^2 \text{ s.t. } H \in [-1, 1]^{n \times q}.$$

- **Algorithm:** random coordinate descent using Newton directions
  1. Randomly select an entry  $H_{ij}$  in  $H$  to update while keeping other entries fixed
  2. Approximate the objective function by second-order Taylor expansion w.r.t.  $H_{ij}$   
Calculate a step-size  $d$  and update  $H_{ij}$  by  $H_{ij} \leftarrow H_{ij} + d$
  3. Repeat 1 and 2 until stopping criterion is satisfied

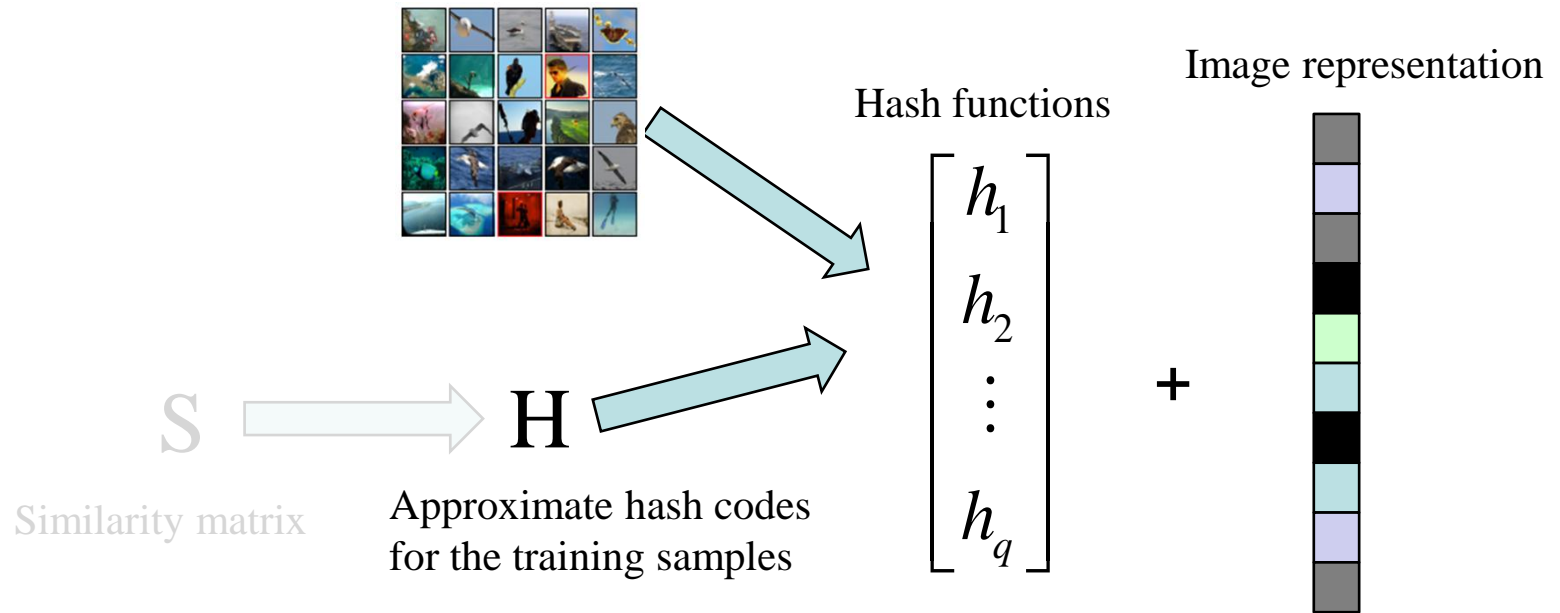
**The time complexity of the whole algorithm is  $O(Tqn^2)$  with small  $T$ ,  $q$ .**

$n$ : the number of training samples

$q$ : hash code length (less than 64 in our experiments)

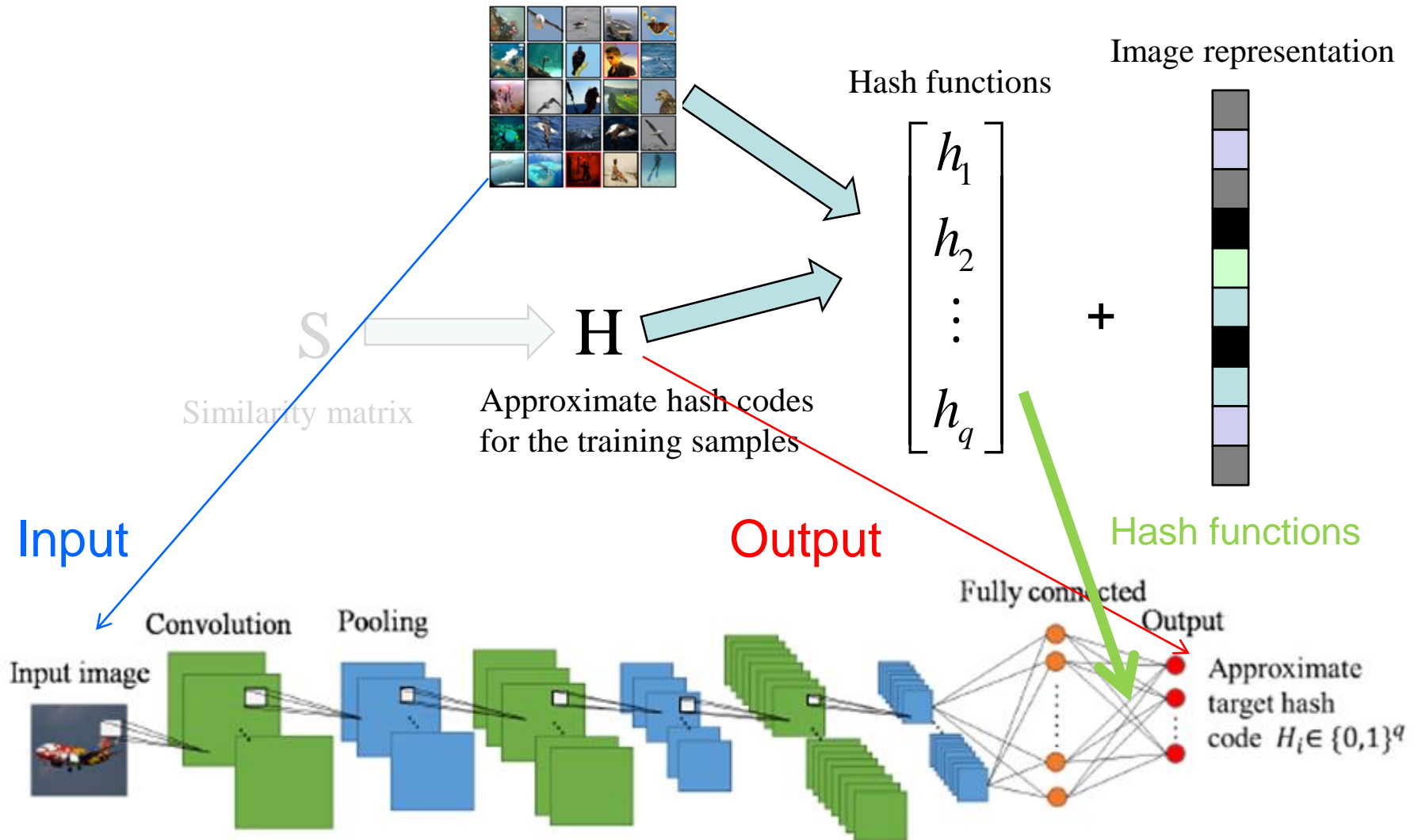
$T$ : iterations (less than 5 in our experiments)

# Stage 2: Learning Hash Functions

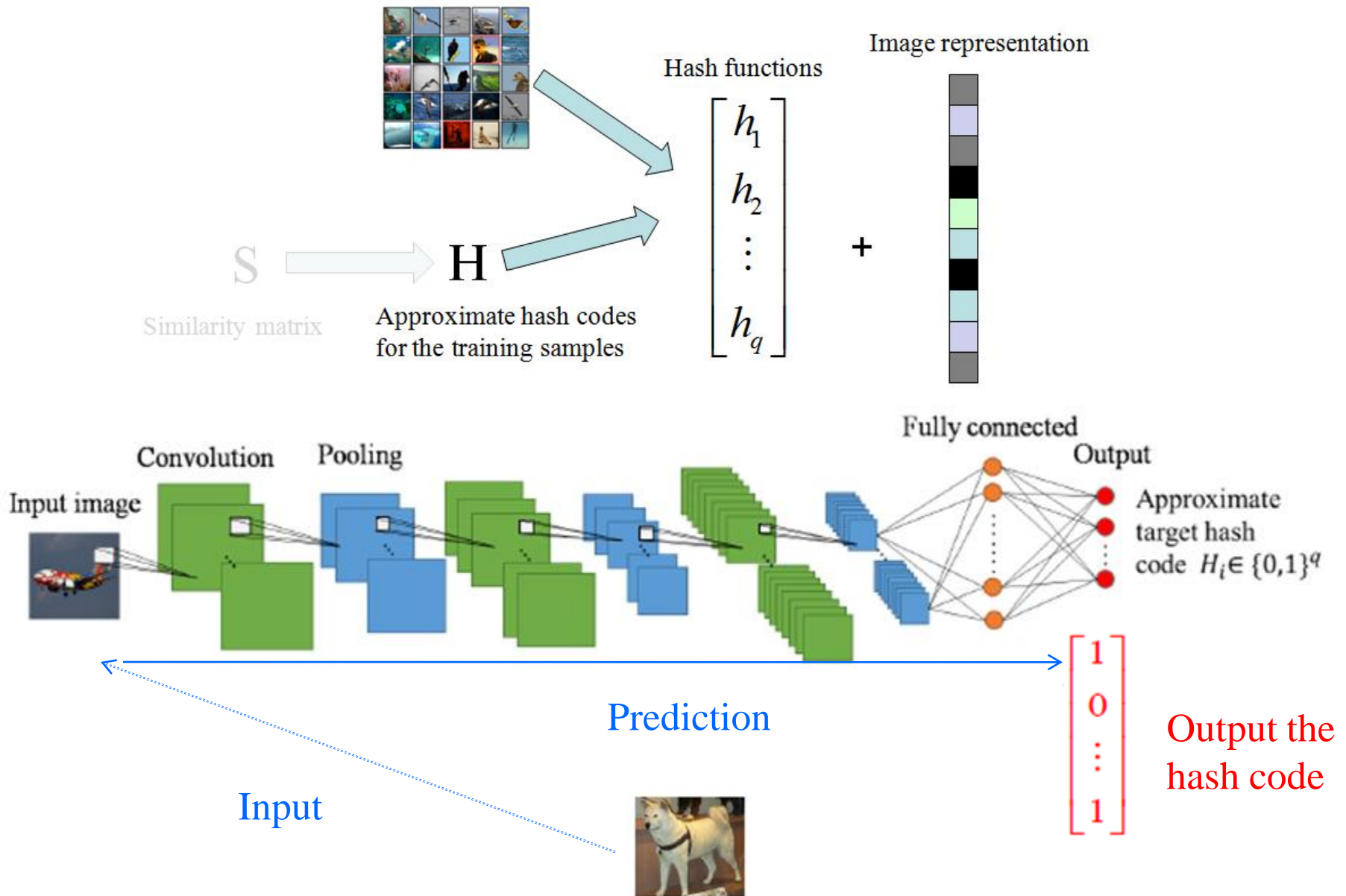


- It is a multi-label binary classification problem that is solved by deep convolutional neural networks.
- It learns hash functions as well as image features.
- We propose two methods: CNNH and CNNH+.

# Method 1: CNNH

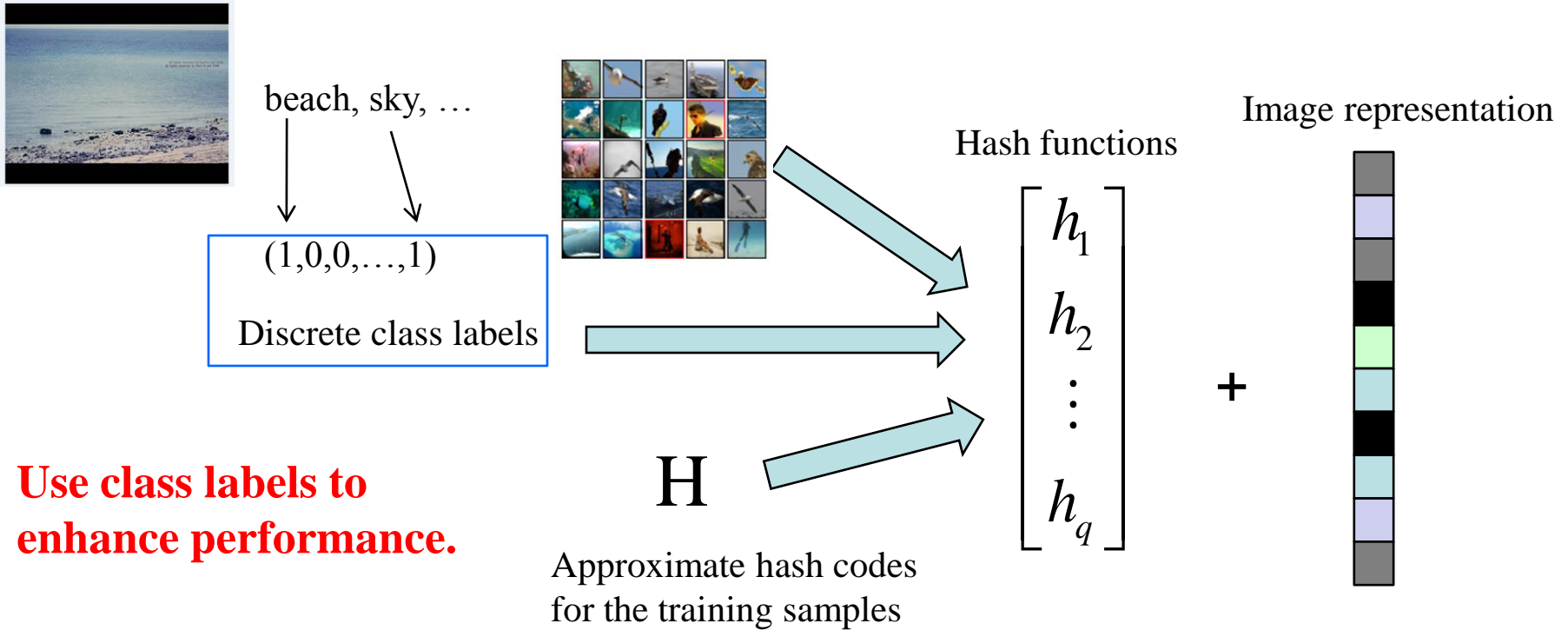


# Method 1: CNNH

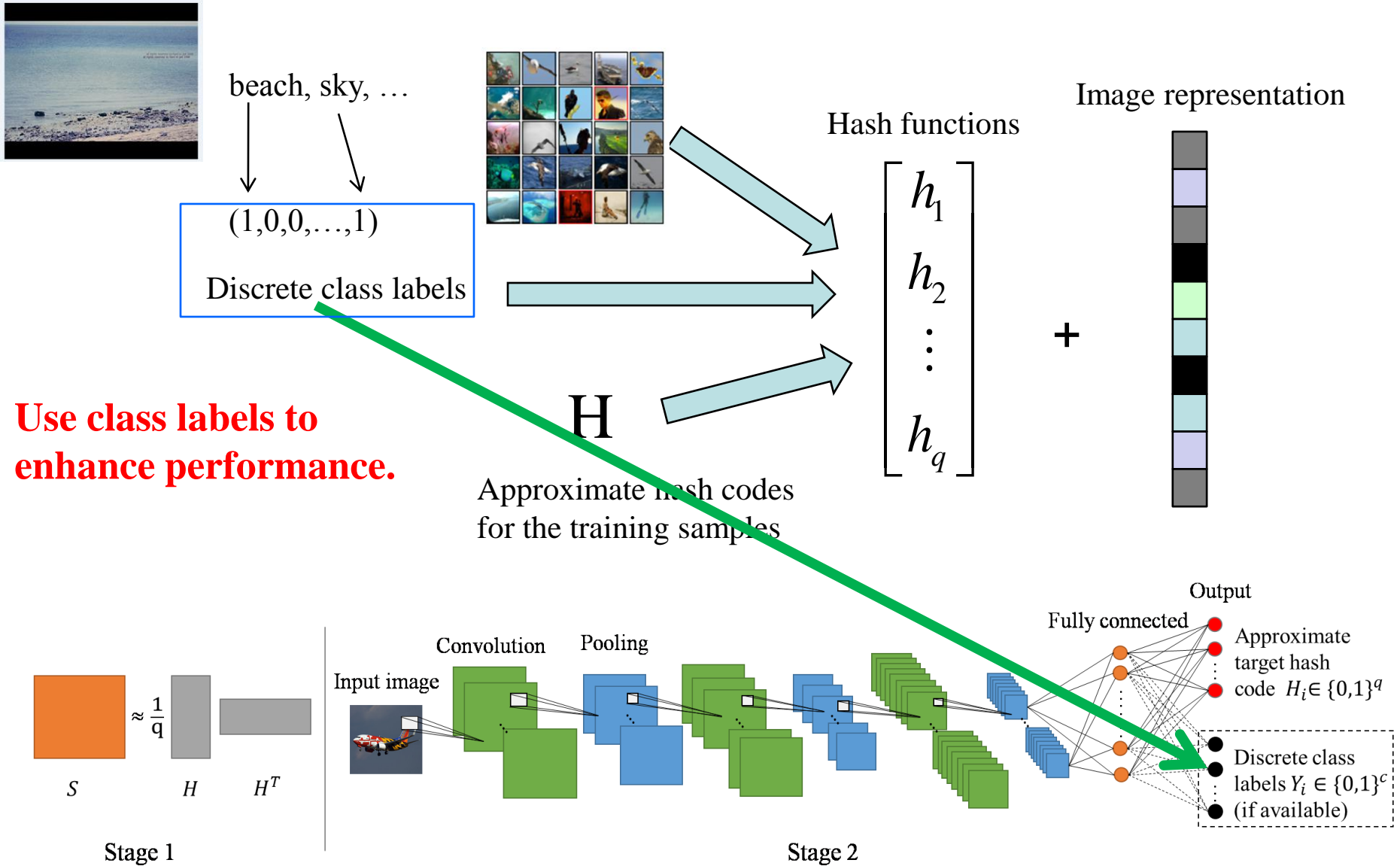




# Method2: CNNH+

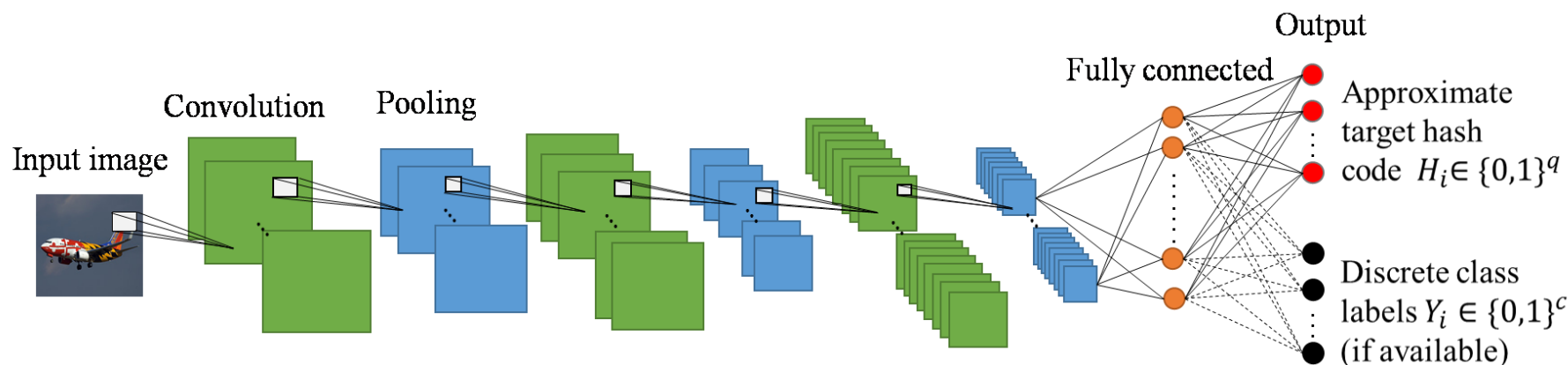


# Method2: CNNH+

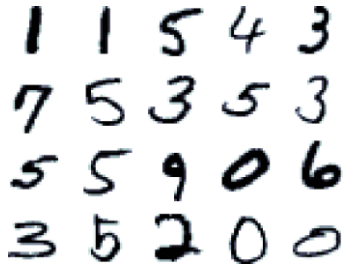


# Details of the Deep Convolutional Networks

- We adopt the architecture of [Krizhevsky, NIPS 2012] as our basic framework.
- Our network has three convolutional-pooling layers with rectified linear activation, max pooling and local contrast normalization, a standard fully connected layer, and an output layer with softmax activation.
- We use 32, 64, 128 filter (with the size  $5 \times 5$ ) in the 1<sup>st</sup>, 2<sup>nd</sup> and 3<sup>rd</sup> convolutional layer, respectively.
- We use dropout with a rate of 0.5.



# Datasets



**MNIST**: 70,000 greyscale images (in size 28\*28) of handwritten digits from '0' to '9'



**CIFAR10**: 60,000 color tinny images (in size 32\*32) that are categorized in 10 classes



**NUS-WIDE**: about 270,000 images collected from the web. It is a multi-label dataset.

# Baseline Methods

## Unsupervised methods

LSH [Gionis et al. VLDB,1999]

SH [Weiss and Torralba. NIPS,2008]

ITQ [Gong and Lazebnik. CVPR,2011]

MLH [Norouzi and Blei. ICML,2011]

BRE [Kulis and Darrell. NIPS,2009]

ITQ-CCA [Gong and Lazebnik. CVPR,2011]

KSH [Liu et al. CVPR,2012]

## Supervised methods

# Evaluation Metrics

Precision:

$$precision = \frac{\#\{\text{retrieved relevant images}\}}{\#\{\text{retrieved images}\}}$$

Recall:

$$recall = \frac{\#\{\text{retrieved relevant images}\}}{\#\{\text{all relevant images}\}}$$

Mean Average Precision (MAP):

$$MAP = \frac{1}{q} \sum_i AP_i$$

$$AP = \frac{\sum_n P @ n \times I\{\text{image } n \text{ is relevant}\}}{\#\{\text{retrieved relevant image}\}}$$

$$P @ n = \frac{\#\{\text{relevant images in top } n \text{ result}\}}{n}$$

# Experimental Results

MAP of Hamming ranking on **MNIST** w.r.t. different number of bits

Methods	12-bit	24-bit	32-bit	48-bit
<b>CNNH+</b>	0.969	0.975	0.971	0.975
<b>CNNH</b>	0.957	0.963	0.956	0.960
<b>KSH</b>	0.872	0.891	0.897	0.900
<b>ITQ-CCA</b>	0.659	0.694	0.714	0.726
<b>MLH</b>	0.472	0.666	0.652	0.654
<b>BRE</b>	0.515	0.593	0.613	0.634
<b>SH</b>	0.265	0.267	0.259	0.250
<b>ITQ</b>	0.388	0.436	0.422	0.429
<b>LSH</b>	0.187	0.209	0.235	0.243

relative increase of 8.2%~11.1%

# Experimental Results

MAP of Hamming ranking on **CIFAR10** w.r.t. different number of bits

Methods	12-bit	24-bit	32-bit	48-bit
<b>CNNH+</b>	0.465	0.521	0.521	0.532
<b>CNNH</b>	0.439	0.511	0.509	0.522
<b>KSH</b>	0.303	0.337	0.346	0.356
<b>ITQ-CCA</b>	0.264	0.282	0.288	0.295
<b>MLH</b>	0.182	0.195	0.207	0.211
<b>BRE</b>	0.159	0.181	0.193	0.196
<b>SH</b>	0.131	0.135	0.133	0.130
<b>ITQ</b>	0.162	0.169	0.172	0.175
<b>LSH</b>	0.121	0.126	0.120	0.120

relative increase of 49.4%~54.6%



# Experimental Results

MAP of Hamming ranking on **NUSWIDE** w.r.t. different number of bits

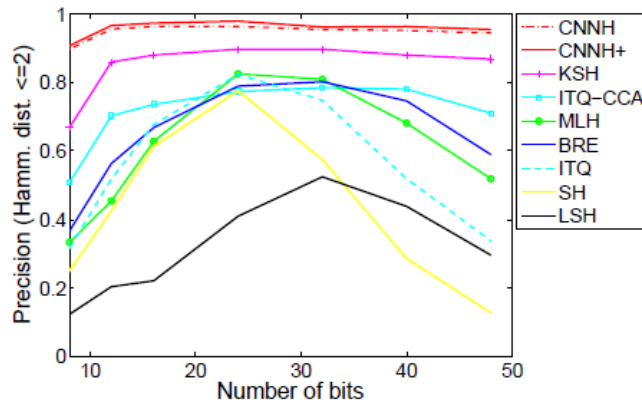
Methods	12-bit	24-bit	32-bit	48-bit
<b>CNNH+</b>	0.623	0.630	0.629	0.625
<b>CNNH</b>	0.611	0.618	0.625	0.608
<b>KSH</b>	0.556	0.572	0.581	0.588
<b>ITQ-CCA</b>	0.435	0.435	0.435	0.435
<b>MLH</b>	0.500	0.514	0.520	0.522
<b>BRE</b>	0.485	0.525	0.530	0.544
<b>SH</b>	0.433	0.426	0.426	0.423
<b>ITQ</b>	0.452	0.468	0.472	0.477
<b>LSH</b>	0.403	0.421	0.426	0.441

relative increase of 6.3%~12.1%

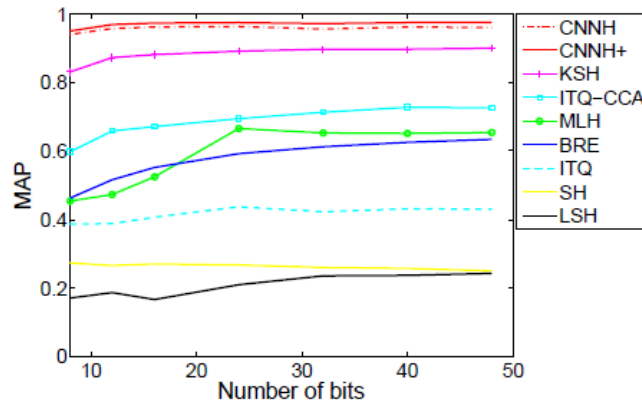
# Experimental Results

## Results on MNIST

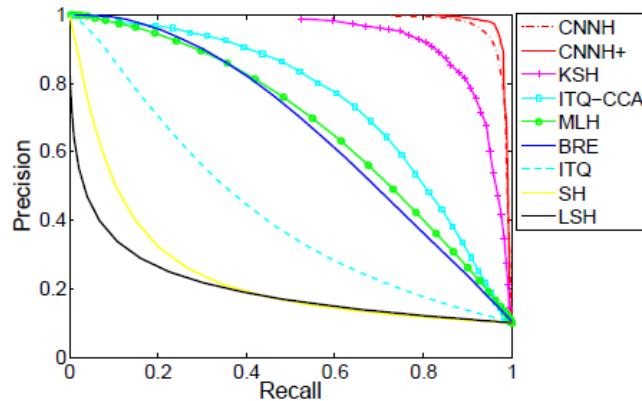
- (a) precision within curves Hamming radius 2 (b) MAP curves within Hamming radius 2  
(c) precision-recall curves with 48 bits (d) precision curves with 48 bits



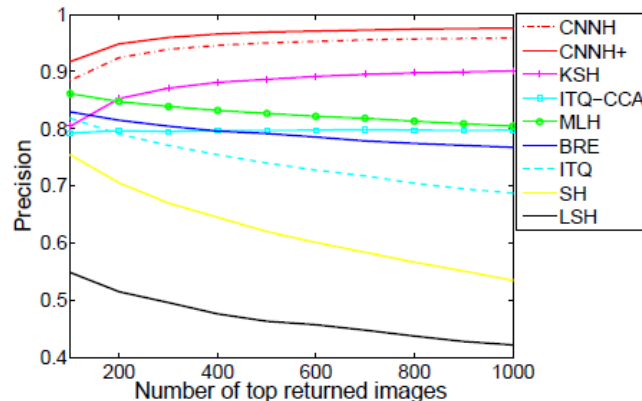
(a)



(b)



(c)

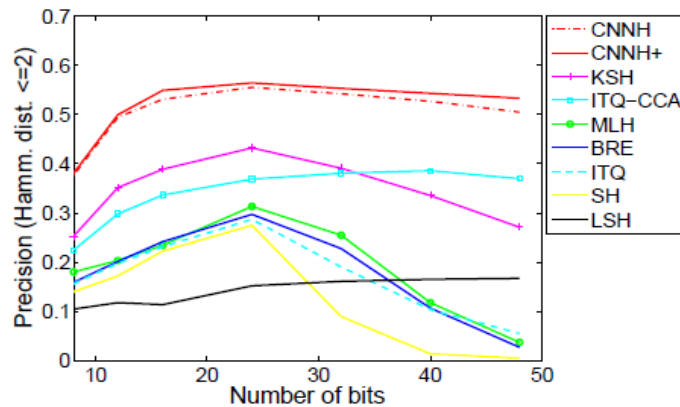


(d)

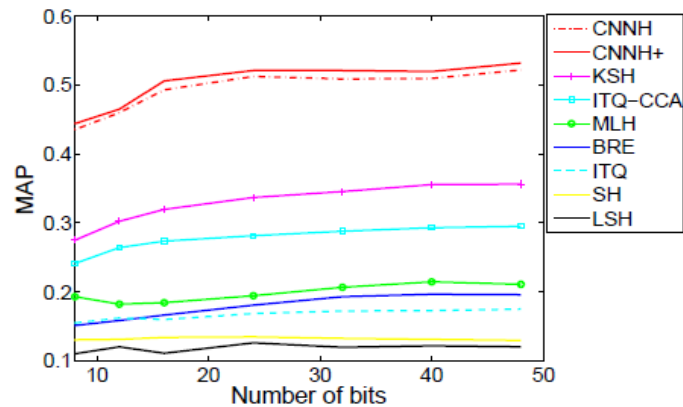
# Experimental Results

## Results on CIFAR-10

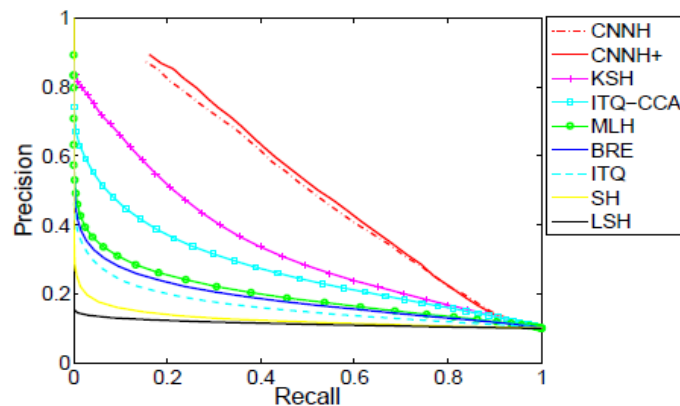
- (a) precision curves within Hamming radius 2 (b) MAP curves within Hamming radius 2  
(c) precision-recall curves with 48 bits (d) precision curves with 48 bits



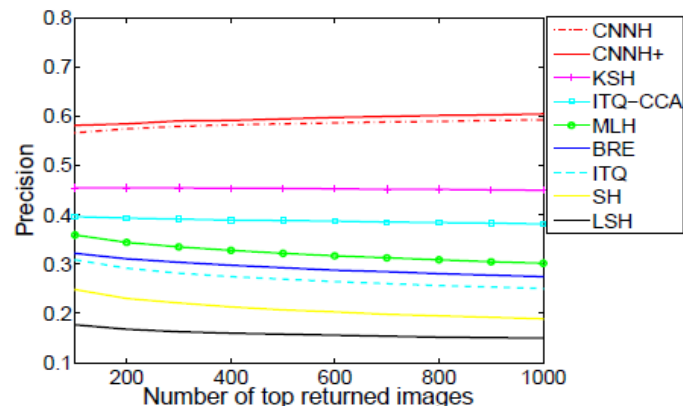
(a)



(b)



(c)

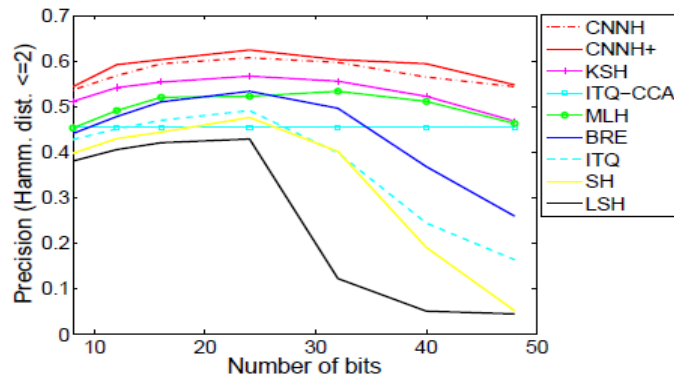


(d)

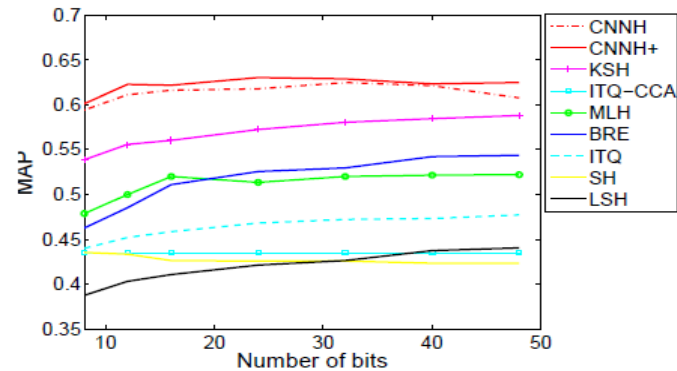
# Experimental Results

## Results on NUS-WIDE

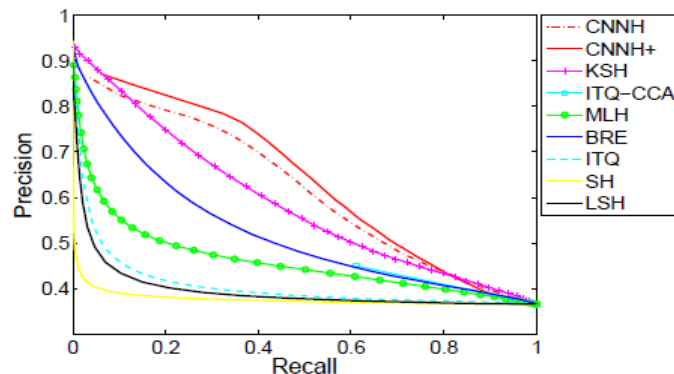
- (a) precision curves within Hamming radius 2 (b) MAP curves within Hamming radius 2  
(c) precision-recall curves with 48 bits (d) precision curves with 48 bits



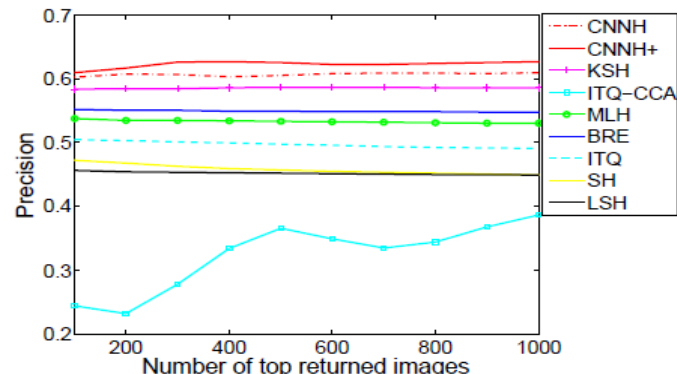
(a)



(b)



(c)

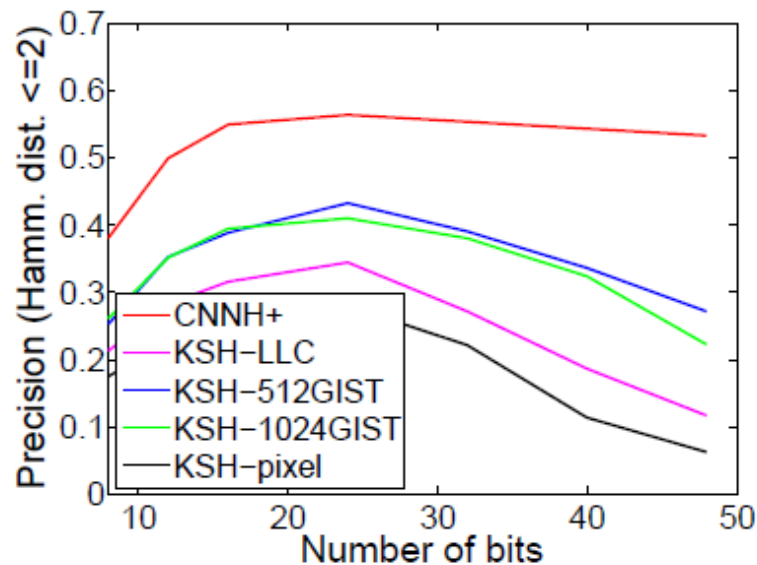


(d)

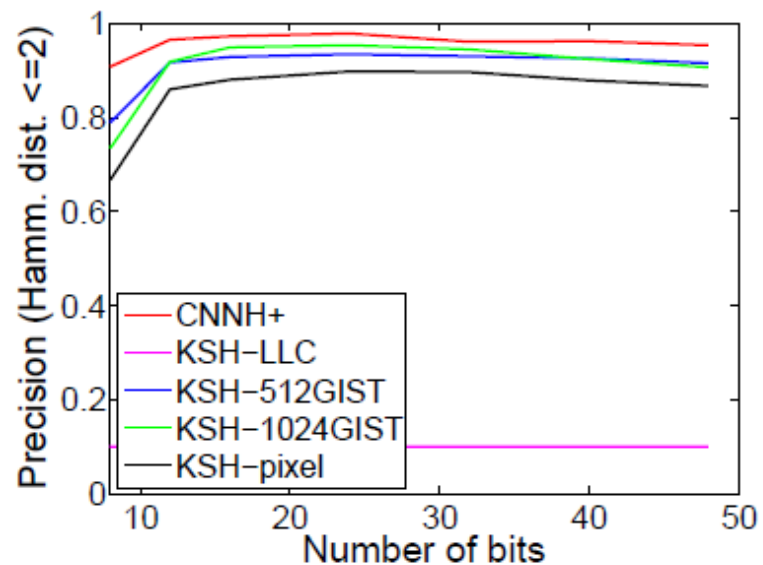
# Experimental Results

CNNH+ vs. KSH with different hand-crafted features

(a) Results on CIFAR-10    (b) Results on MNIST



(a)



(b)

The performances of KSH with different features are inferior to those of CNNH+.

**Thank you!**