# Similarity-Based Visual Search

**Deep Hash Neural Network Implementation** based on

Hanjiang Lai, Yan Pan, Ye Liu, and Shuicheng Yan. Simultaneous feature learning and hash coding with deep neural networks, CVPR 2015.

Yupan Huang (**Sun Yat-SenUniversity**) | 2017.4.14

# Finding Similar Images

Assume $I$ to be the image space. The goal of hash learning for images is to learn a mapping $F: I \rightarrow \{0,1\}^q$, such that an input image $I$ can be encoded into a $q$-bit binary code $F(I)$, with the similarities of images being preserved.



Google Image Search Results

# Motivation

- Extracting informative image features
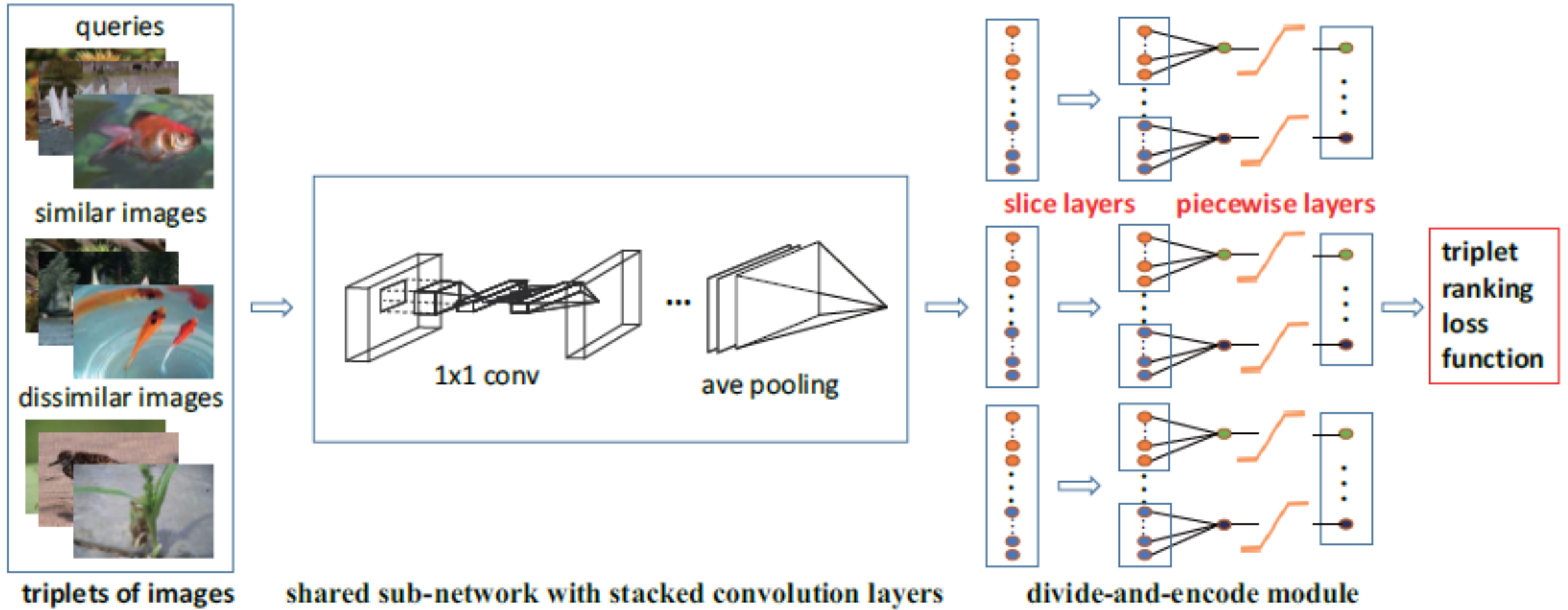  Deep Convolutional Neural Network:
  - Extract discriminative features depend on data
  - Build end-to-end relation between the raw image data and the binary hashing codes

- Learning effective approximate functions
  Hash Code:
  - **Space-saving :**high-dimensional features to lower dimensional space, compact bitwise representation
  - **Speedup :**binary pattern matching or Hamming distance measurement

# Approach



Three building blocks of DNNH

# Part 1 Triplet Ranking Loss and Optimization

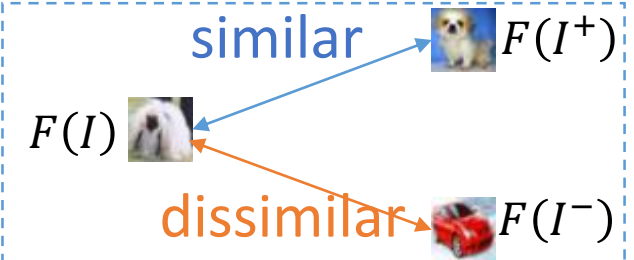$$\hat{l}_{triplet}\left(F(I), F(I^+), F(I^-)\right)$$
$$= \max\left(0, 1 - \left(\left\|F(I) - F(I^-)\right\|_H - \left\|F(I) - F(I^+)\right\|_H\right)\right)$$
$$s.t.\ F(I), F(I^+) \in \{0,1\}^q$$

Characterize that one image is more similar to the second image than to the third one

Relaxation

similar  $F(I^+)$

$F(I)$

dissimilar  $F(I^-)$

Forward inference: $f_w(x)$



1x1 conv ... ave pooling

Backward learning: $\nabla f_w(x)$

$$\hat{l}_{triplet}\left(F(I), F(I^+), F(I^-)\right)$$
$$= \max\left(0, \left\|F(I) - F(I^+)\right\|_2^2 - \left\|F(I) - F(I^-)\right\|_2^2 + 1\right)$$
$$s.t.\ F(I), F(I^+) \in [0,1]^q$$

$$\frac{\partial l}{\partial b} = (2b^- - 2b^+) \times I_{\left\|b-b^+\right\|_2^2 - \left\|b-b^-\right\|_2^2 + 1 > 0}$$
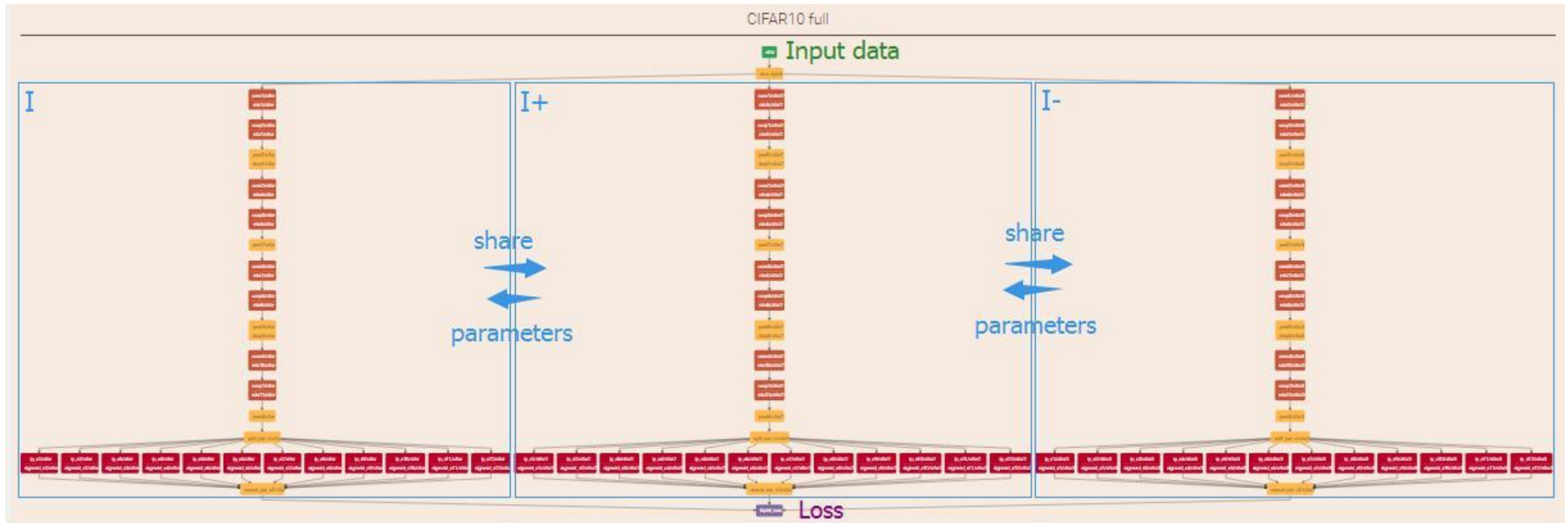
$$\frac{\partial l}{\partial b} = (2b^+ - 2b) \times I_{\left\|b-b^+\right\|_2^2 - \left\|b-b^-\right\|_2^2 + 1 > 0}$$

$$\frac{\partial l}{\partial b} = (2b^- - 2b) \times I_{\left\|b-b^+\right\|_2^2 - \left\|b-b^-\right\|_2^2 + 1 > 0}$$
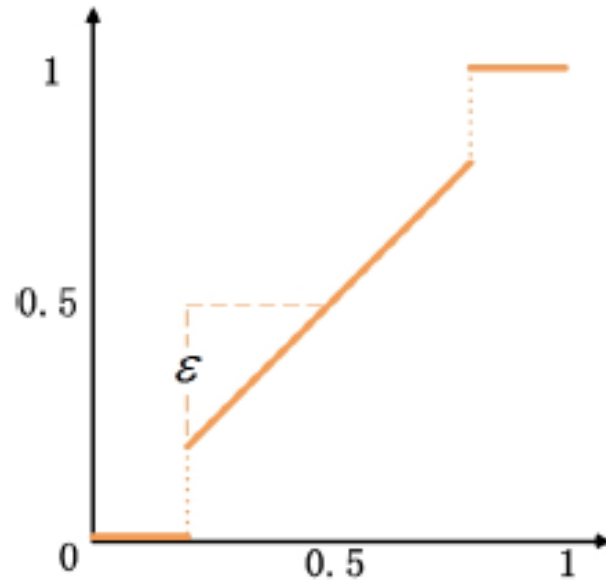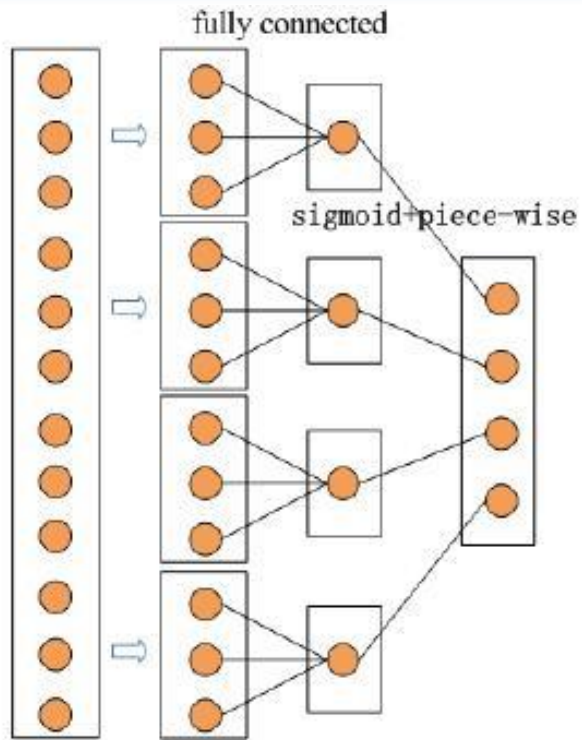
5

# Part 2 Shared SubNetwork with Stacked Convolution Layers

- Parameter sharing can significantly reduce the number of parameters in the whole architecture.

http://ethereon.github.io/netscope/#/gist/2f44d589dbe4355ad2f9374344527b6b

# Part 3 Divide and Encode Module

fully connected

sigmoid+piece-wise

- Separated slice of features: Reduce the redundancy among the hash bits

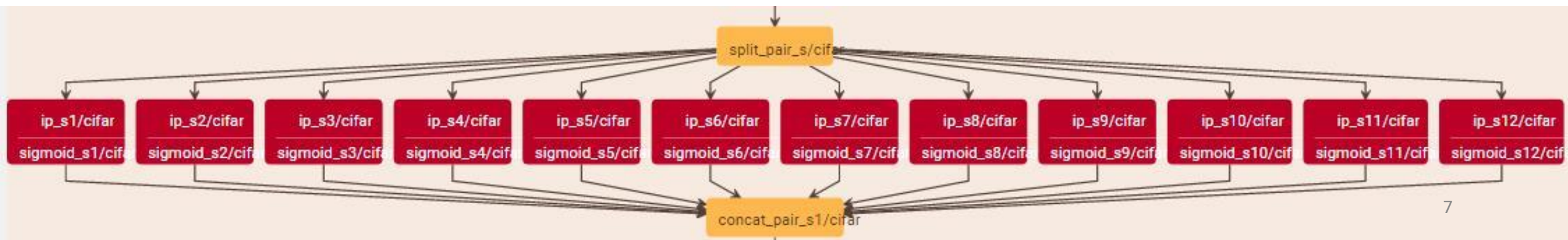- Piece-wise threshold: encourage the output of binary hash bits:

split_pair_s/cifar

| ip_s1/cifar | ip_s2/cifar | ip_s3/cifar | ip_s4/cifar | ip_s5/cifar | ip_s6/cifar | ip_s7/cifar | ip_s8/cifar | ip_s9/cifar | ip_s10/cifar | ip_s11/cifar | ip_s12/cifar |
| sigmoid_s1/cifa | sigmoid_s2/cif | sigmoid_s3/cif | sigmoid_s4/cif | sigmoid_s5/cif | sigmoid_s6/cif | sigmoid_s7/cif | sigmoid_s8/cif | sigmoid_s9/cif | sigmoid_s10/cif | sigmoid_s11/cif | sigmoid_s12/cif |

concat_pair_s1/cifar

7

# Image Retrieval and Model Evaluation



image  shared sub network  divide-and-encode  quantization
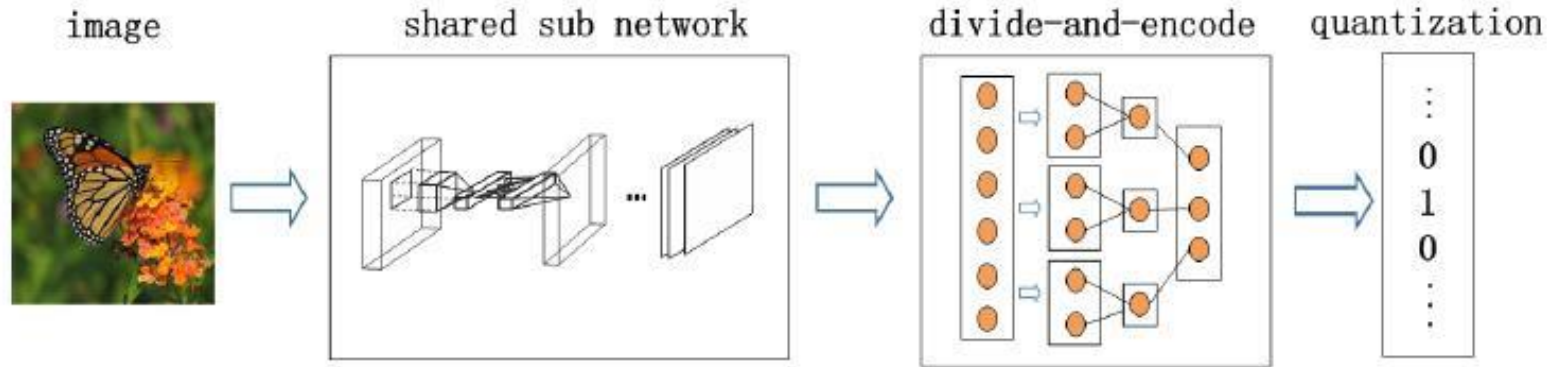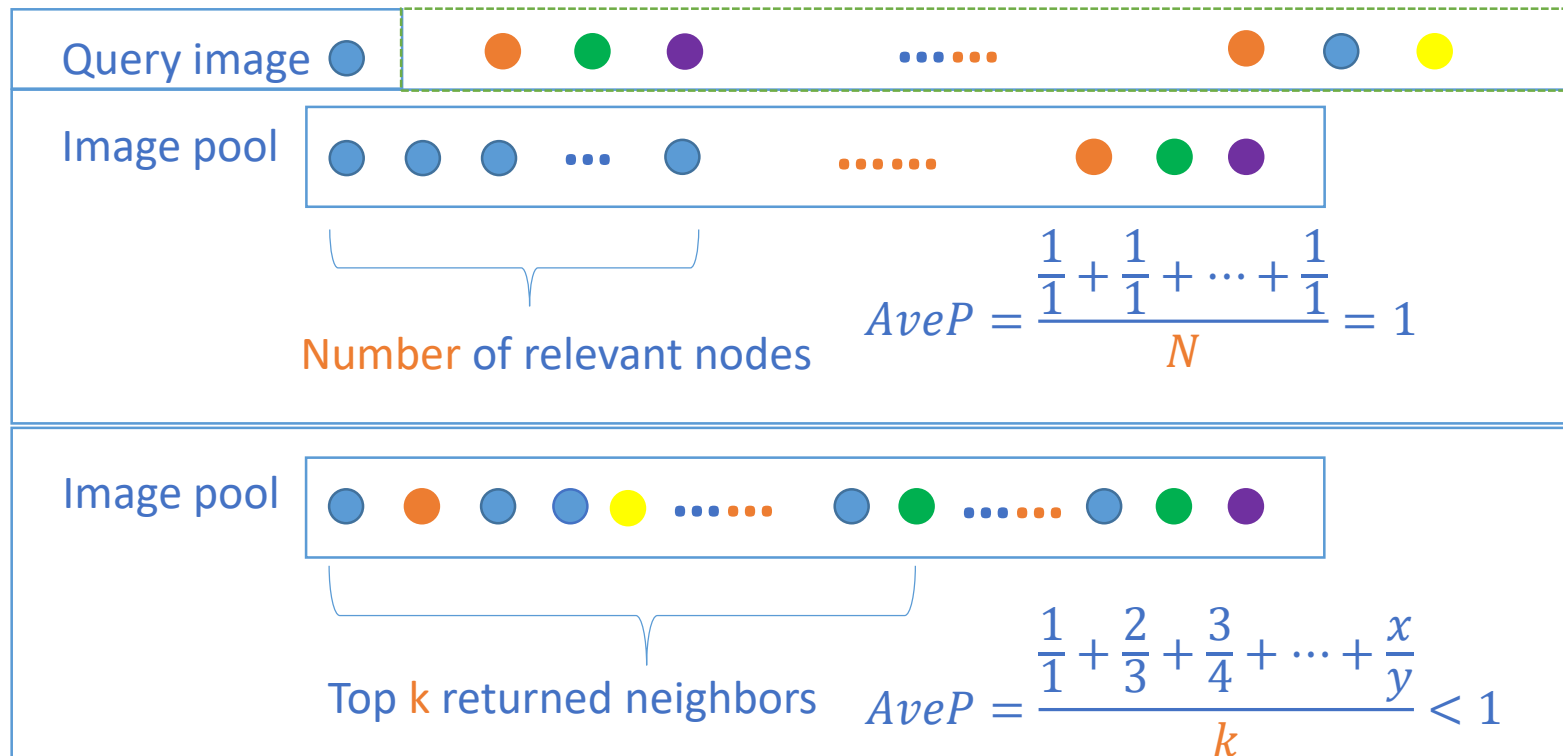
Given a query image, the retrieval list of images is produced by sorting the hamming distances between the query image and images in search pool.

Query image

Image pool

Number of relevant nodes

$$AveP = \frac{\frac{1}{1} + \frac{1}{1} + \cdots + \frac{1}{1}}{N} = 1$$

Image pool

Top k returned neighbors

$$AveP = \frac{\frac{1}{1} + \frac{2}{3} + \frac{3}{4} + \cdots + \frac{x}{y}}{k} < 1$$
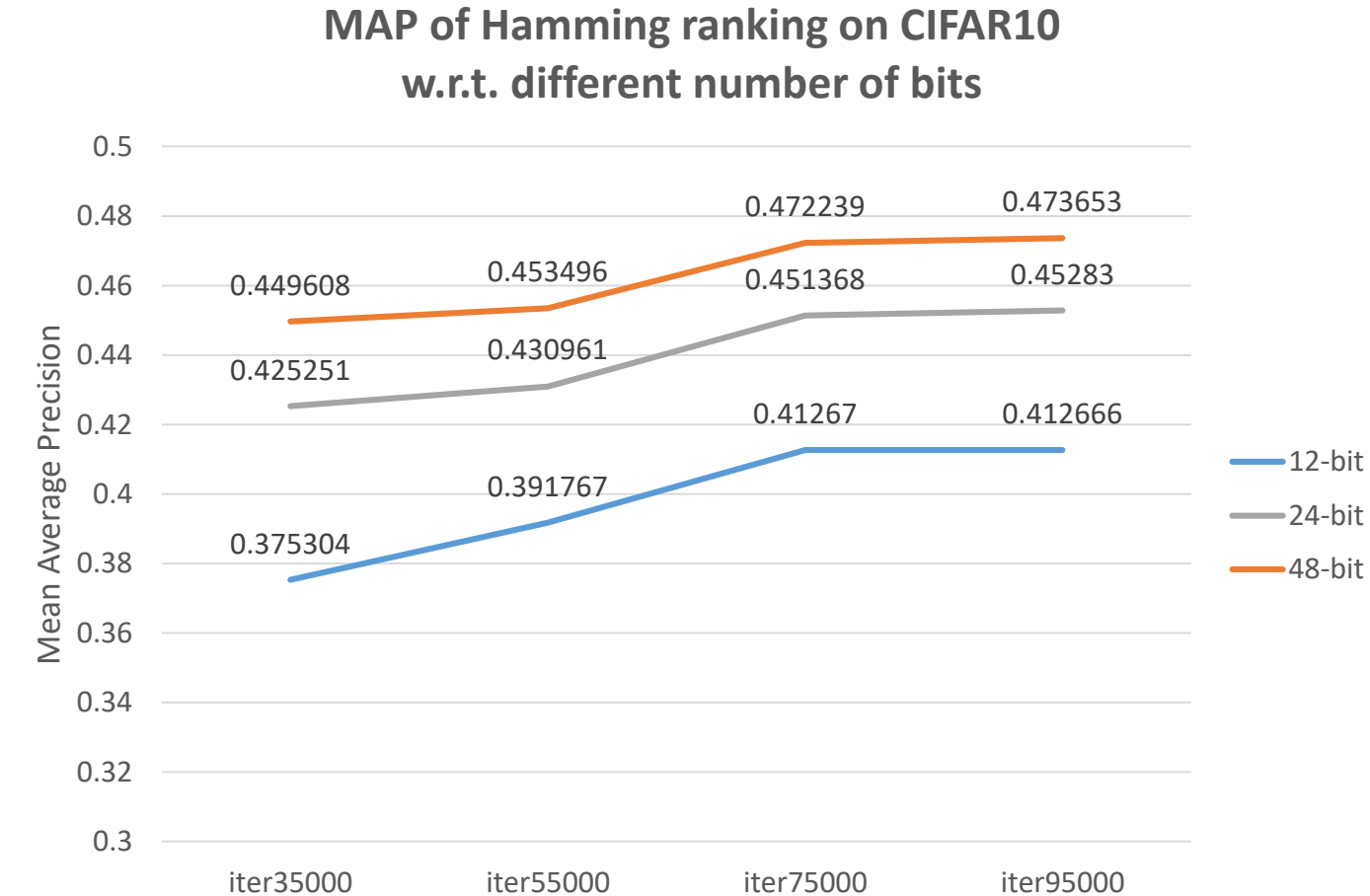
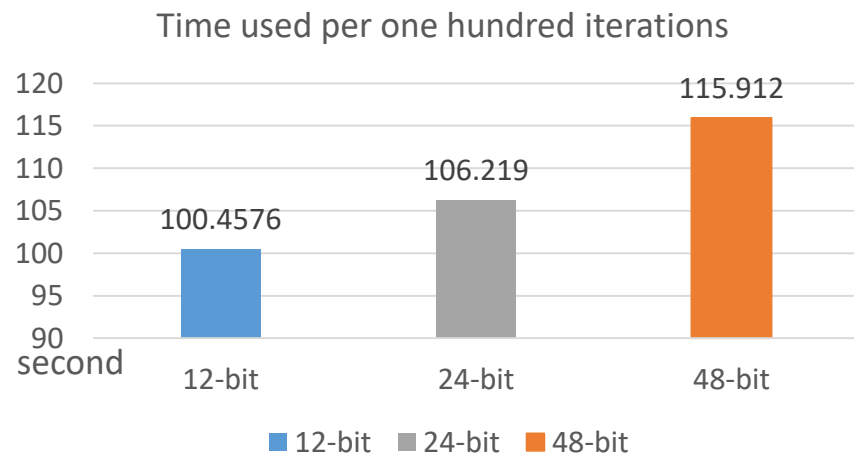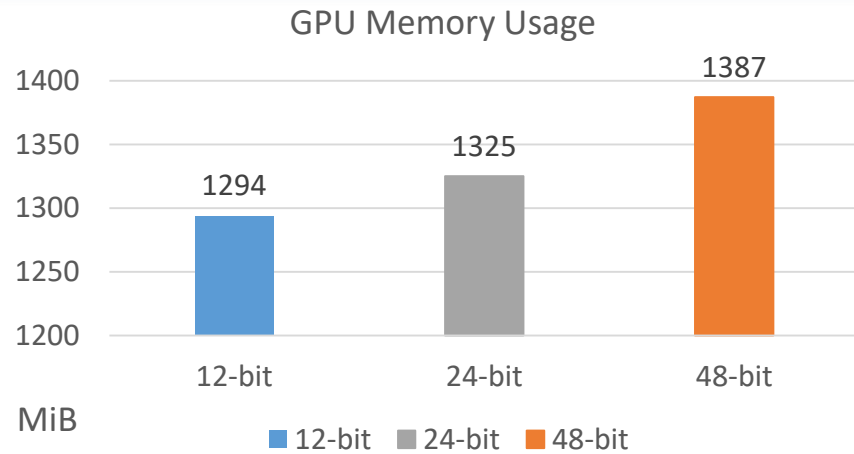$$MAP = \frac{\sum_{q=1}^{Q} AveP(q)}{Q}$$

$$AveP = \frac{\sum_{k=1}^{N}(P(k) \times rel(k))}{\min(N, k)}$$

Mean Average Precision(MAP)

# Implementation on Caffe framwork

- **Deploy**: Given the definition of loss layer, deploy the deep hashing pipeline on linux.

- **Train**: Write prototxt to define DNNH and bash files to execute for training on preprocessed triplet CIFAR-10 dataset.

- **Test**: Write prototxt to encode images into 12-bit, 32-bit and 48-bit seperately and bash files to execute for image retrieval.

- **Evaluate**: Implement the metric of mean average precision (mAP) for evaluation.

# Results and Analysis

## GPU Memory Usage



| bits | MiB |
|------|-----|
| 12-bit | 1294 |
| 24-bit | 1325 |
| 48-bit | 1387 |

MiB

■ 12-bit  ■ 24-bit  ■ 48-bit

## Time used per one hundred iterations



| bits | second |
|------|--------|
| 12-bit | 100.4576 |
| 24-bit | 106.219 |
| 48-bit | 115.912 |

second

■ 12-bit  ■ 24-bit  ■ 48-bit

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



Mean Average Precision

| | iter35000 | iter55000 | iter75000 | iter95000 |
|--------|-----------|-----------|-----------|-----------|
| 48-bit | 0.449608 | 0.453496 | 0.472239 | 0.473653 |
| 24-bit | 0.425251 | 0.430961 | 0.451368 | 0.45283 |
| 12-bit | 0.375304 | 0.391767 | 0.41267 | 0.412666 |

—— 12-bit  —— 24-bit  —— 48-bit

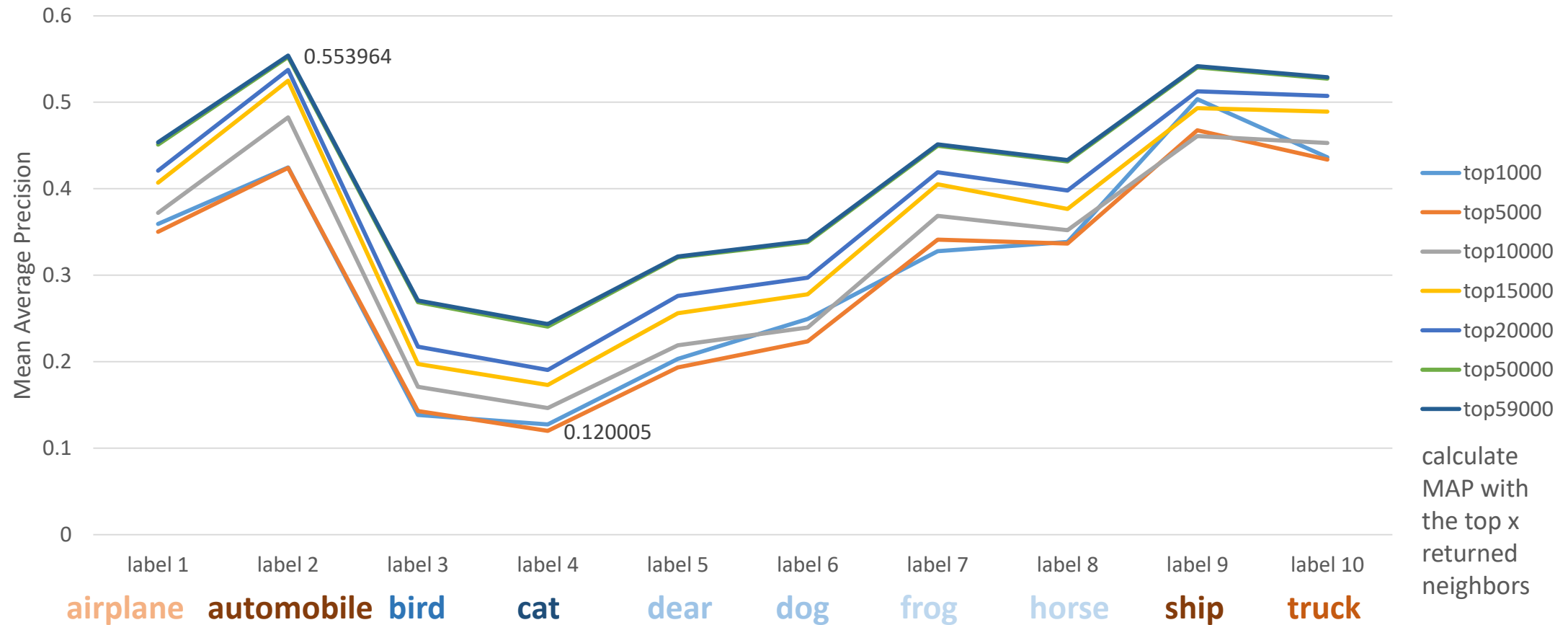Shorter codes have the edge in space and time usage.

Longer codes have higher accuracy.

**128-bit codes**??
Shorter codes may be capable of presenting images and large bits may be overfitting.   Test on Tesla K80 GPU Card
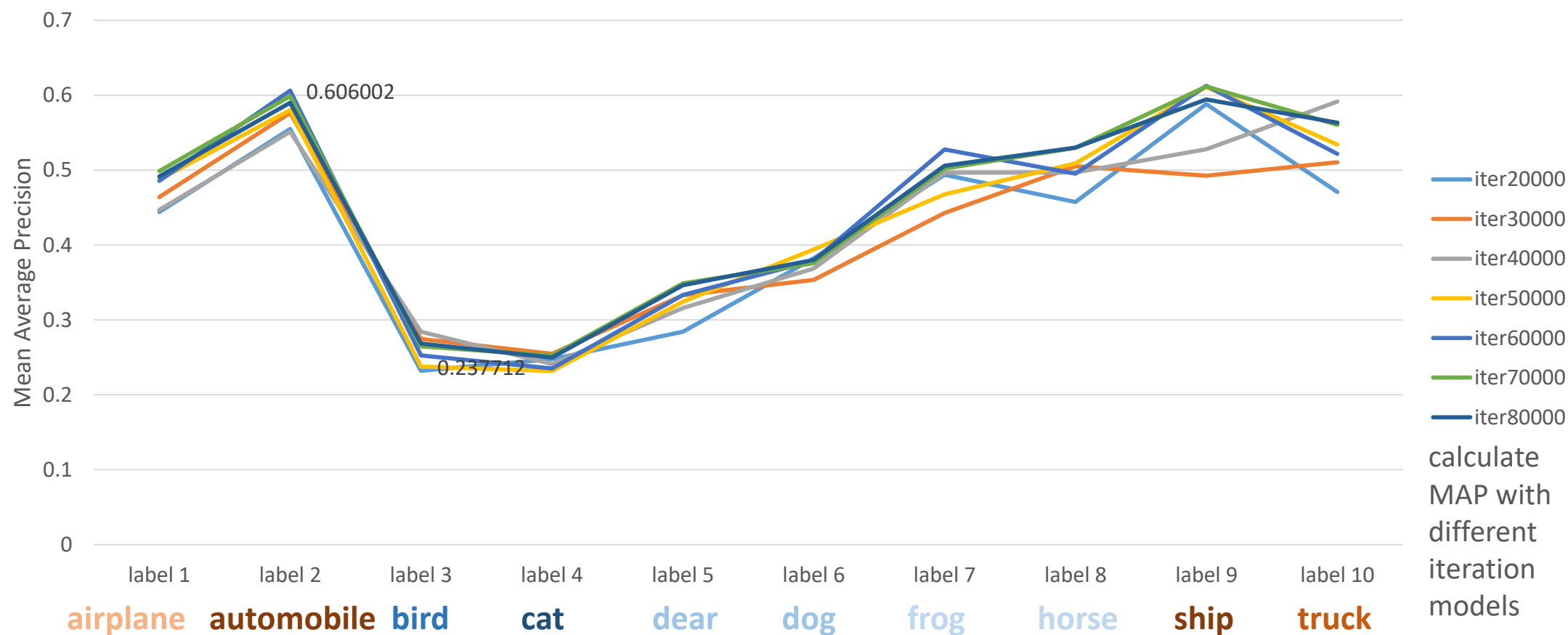
# Results and Analysis

10 labels image's MAP with 12 hash bits, 80000 iterations



Artifacts like automobile, ship, truck, airplane outperform animals like cat, bird, dear dog and so on.

Test on Tesla K80 GPU Card

# Results and Analysis

10 labels image's MAP with 24 hash bits, top 50000 returned neighbors
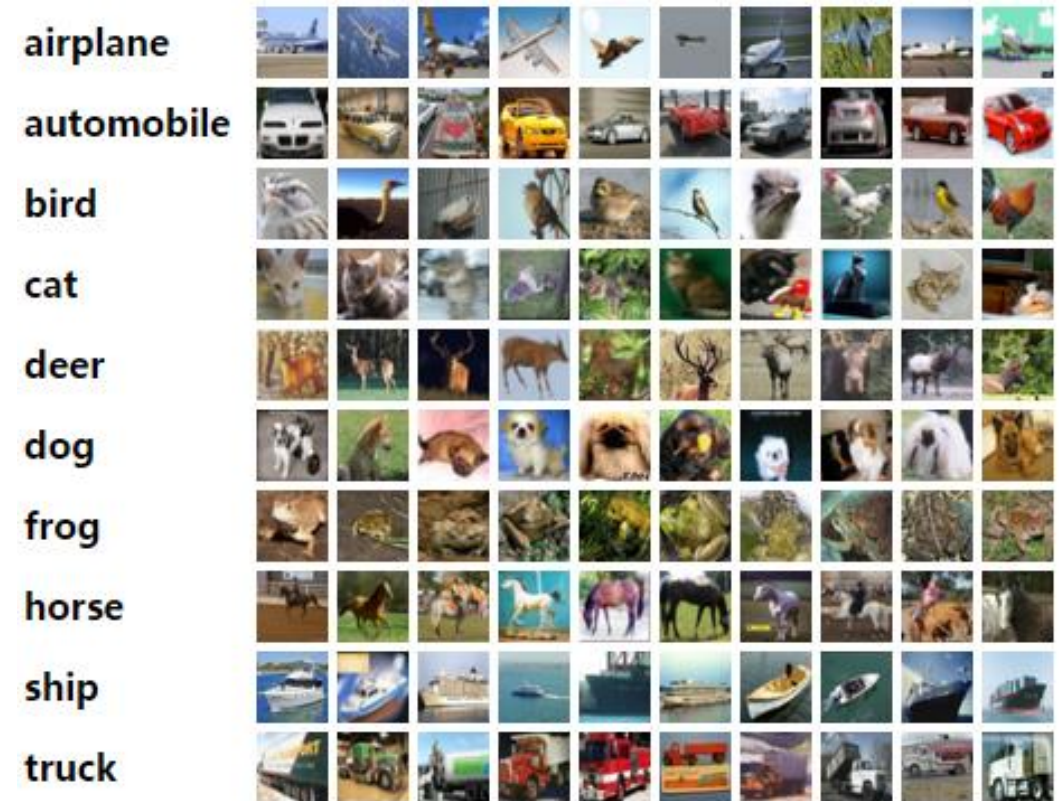


Artifacts like automobile, ship, truck, airplane outperform animals like cat, bird, dear dog and so on.

Test on Tesla K80 GPU Card

# Futher Improvements

- MAP closely relates to labels/classes in the CIFAR-10 dataset.

- Use different margin for different labels/classes.

- *Larger* margin for artifacts.
  *Smaller* margin for animals.



**The CIFAR-10 dataset**

$$\hat{l}_{triplet}\big(F(I), F(I^+), F(I^-)\big)$$
$$= \max\left(0, margin - \left(\big\|F(I) - F(I^-)\big\|_H - \big\|F(I) - F(I^+)\big\|_H\right)\right)$$
$$s.t. \, F(I), F(I^+) \in \{0,1\}^q$$

# Futher Improvements

- Margin design

- Fine-tune on pre-trained deep models for better performance in both speed and precision.

- How about four images a group?
  For example: $\{F(I^{a,b}), F(I^a), F(I^b), F(I^-)\}$
  $F(I^{a,b})$ is image with tag $a$ and $b$, $F(I^a)$ and $F(I^b)$ are images with tag $a$ and $b$ respectively while $F(I^-)$ neither has tag $a$ or $b$.
  Aim to preserve semantic similarity as well.

# Thank you!

For more details:

[Source code](#) in GITHUB.

[Post](#) in blog.