



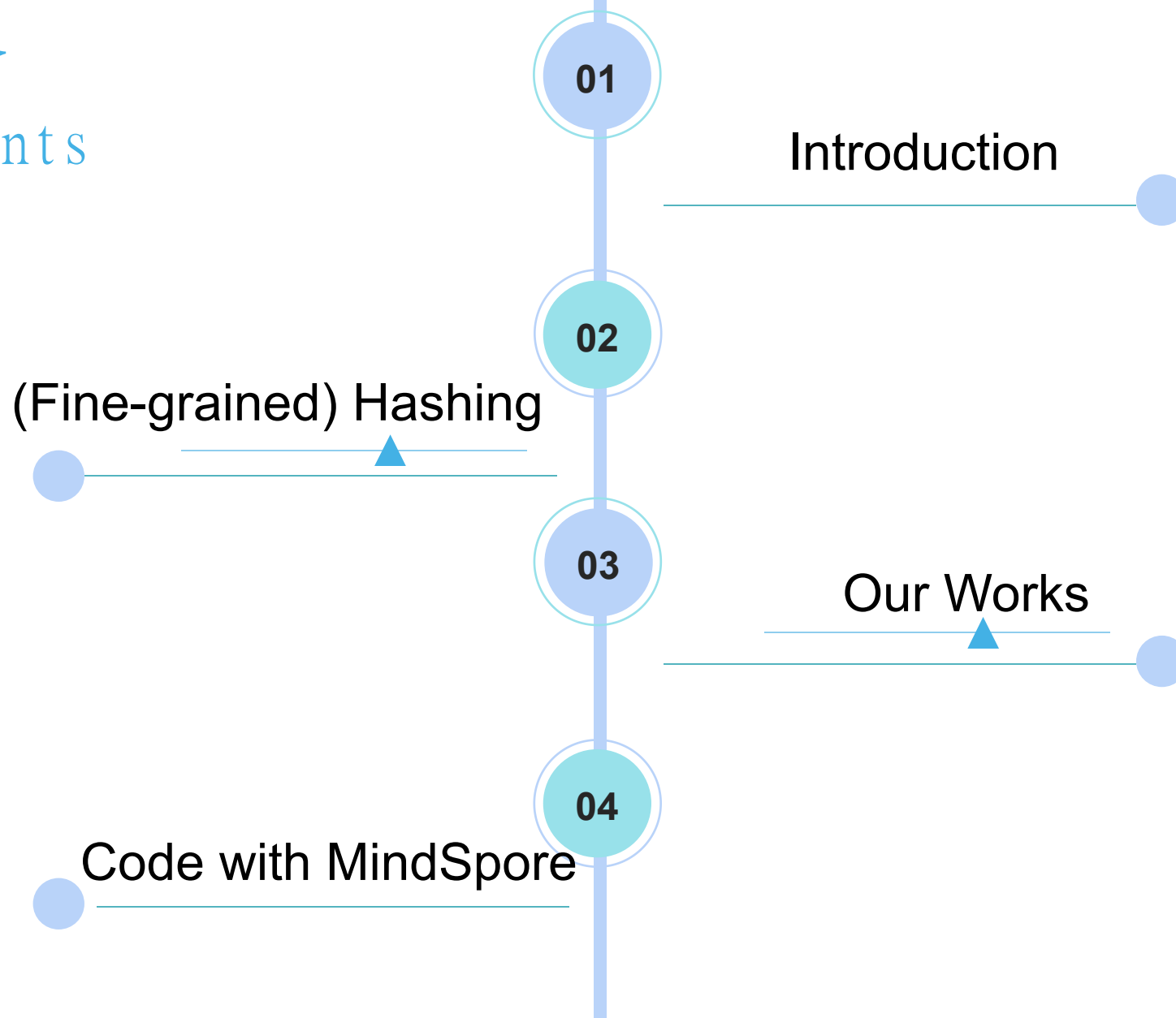
MindSpore

# 基于MindSpore的细 粒度哈希学习

作者：沈阳

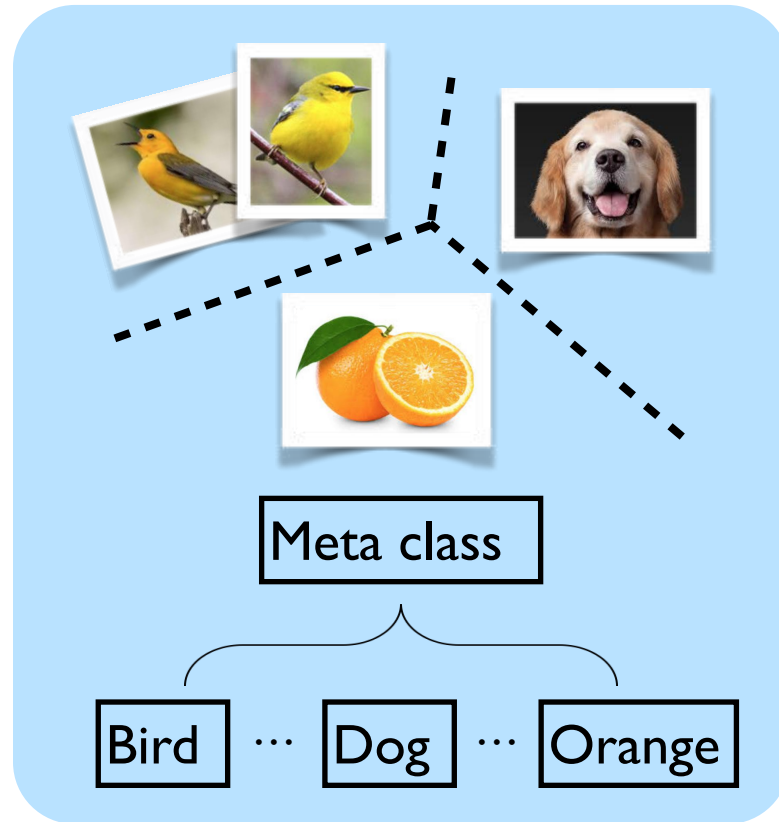
# 目录

contents

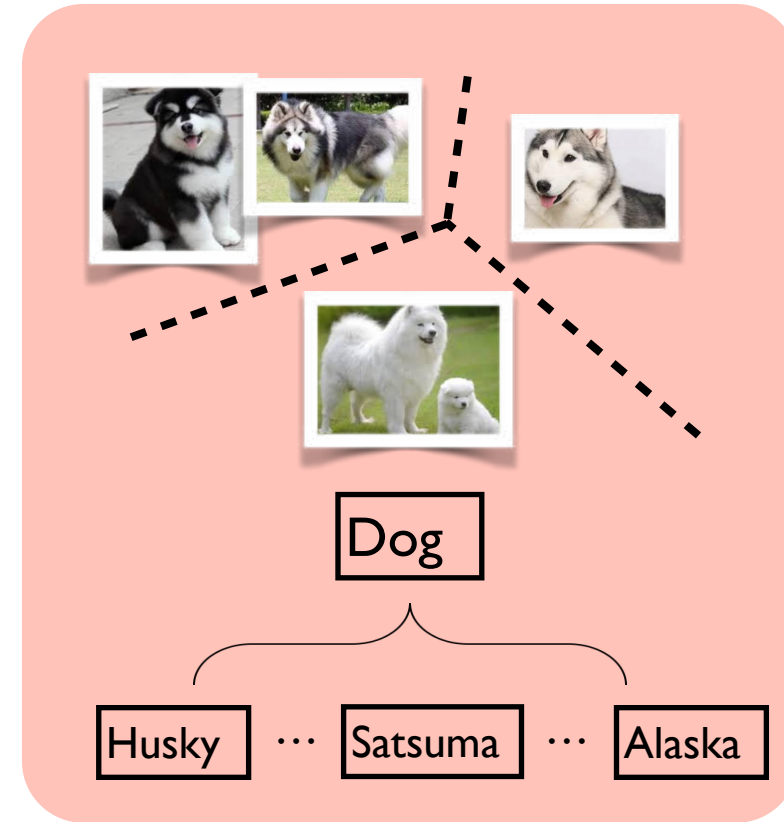


# Introduction

## Fine-grained images vs. generic images



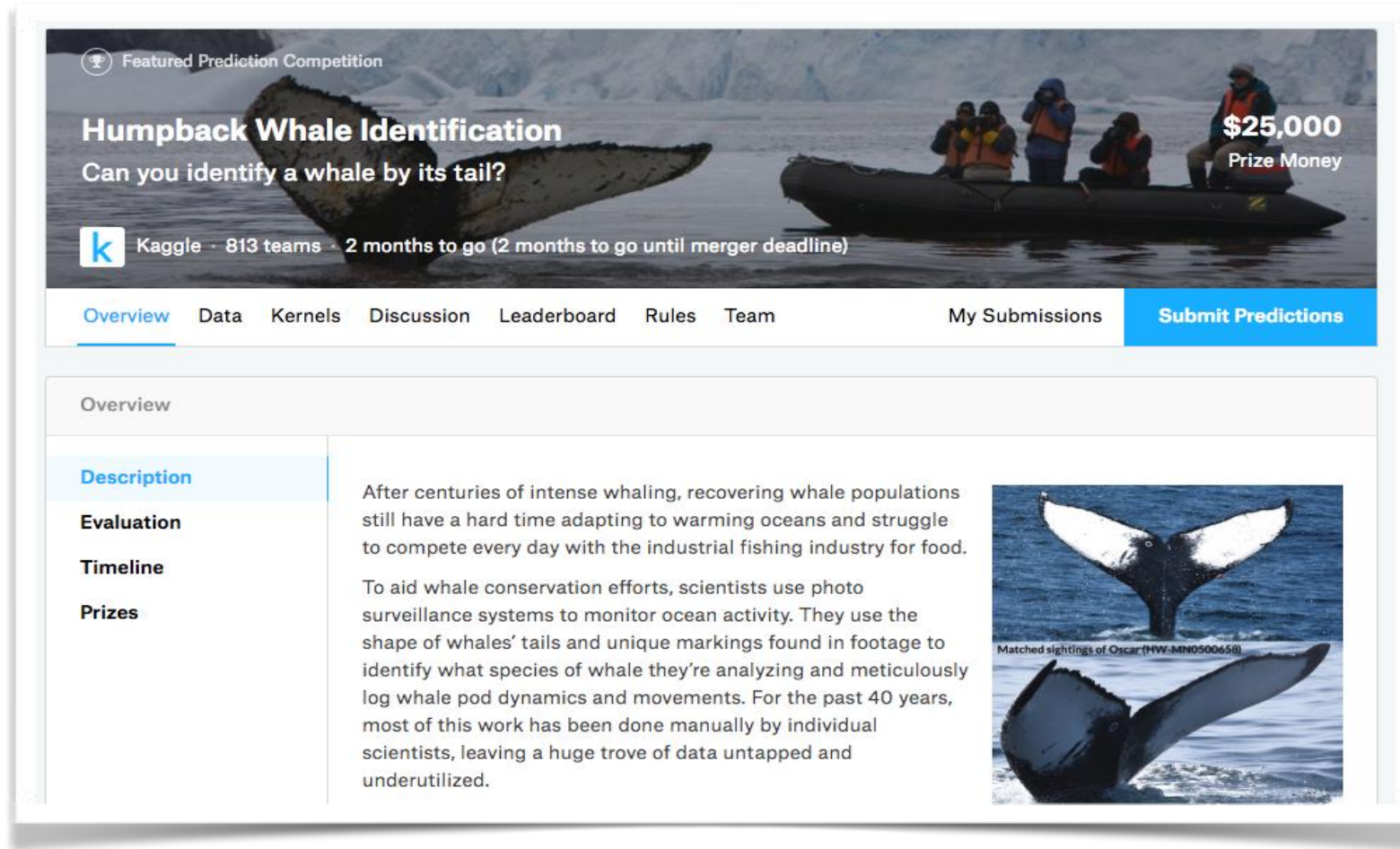
Traditional image tasks  
(Coarse-grained)



Fine-grained image tasks

# Introduction (con't)

## Various real-world applications



The screenshot shows the Kaggle competition page for "Humpback Whale Identification". The header features a large image of a whale's tail and a boat with researchers. Text on the header includes "Featured Prediction Competition", "Humpback Whale Identification", "Can you identify a whale by its tail?", "\$25,000 Prize Money", and "Kaggle · 813 teams · 2 months to go (2 months to go until merger deadline)". Below the header is a navigation bar with links: Overview, Data, Kernels, Discussion, Leaderboard, Rules, Team, My Submissions, and a blue "Submit Predictions" button. The "Overview" section is active, showing a sidebar with "Description", "Evaluation", "Timeline", and "Prizes". The main content area contains text about whale conservation and a photo of two whale tails.

Featured Prediction Competition

## Humpback Whale Identification

Can you identify a whale by its tail?

Kaggle · 813 teams · 2 months to go (2 months to go until merger deadline)

\$25,000 Prize Money

Overview Data Kernels Discussion Leaderboard Rules Team My Submissions **Submit Predictions**

### Overview

**Description**


**Evaluation**

**Timeline**

**Prizes**

After centuries of intense whaling, recovering whale populations still have a hard time adapting to warming oceans and struggle to compete every day with the industrial fishing industry for food.

To aid whale conservation efforts, scientists use photo surveillance systems to monitor ocean activity. They use the shape of whales' tails and unique markings found in footage to identify what species of whale they're analyzing and meticulously log whale pod dynamics and movements. For the past 40 years, most of this work has been done manually by individual scientists, leaving a huge trove of data untapped and underutilized.



Matched sightings of Oscar (HW-MN0500658)



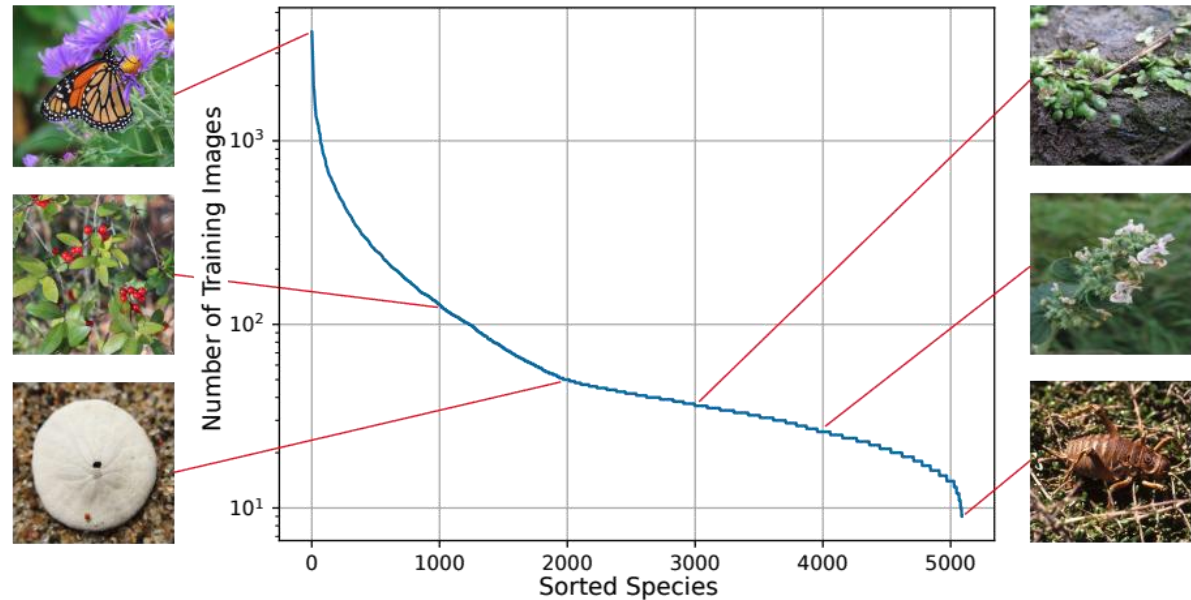
# Introduction (con't)

## Various real-world applications



MindSpore

iNaturalist.org

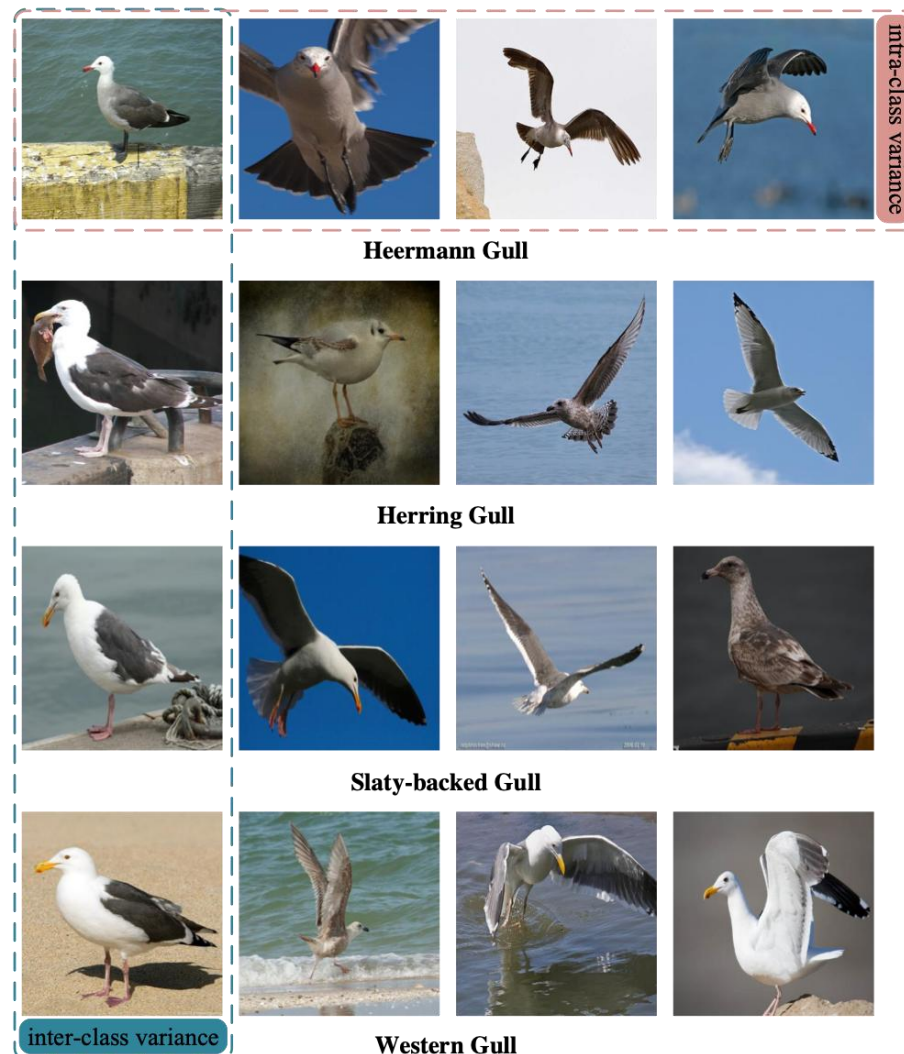


<https://www.kaggle.com/c/inaturalist-challenge-at-fgvc-2017>

# Introduction (con't)

## Challenge of Fine-Grained Image Analysis (FGIA)

**Small inter-class variance**  
**Large intra-class variance**



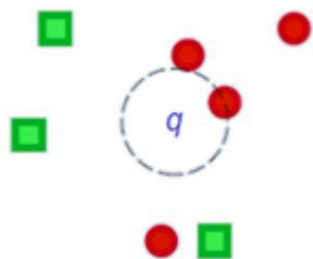
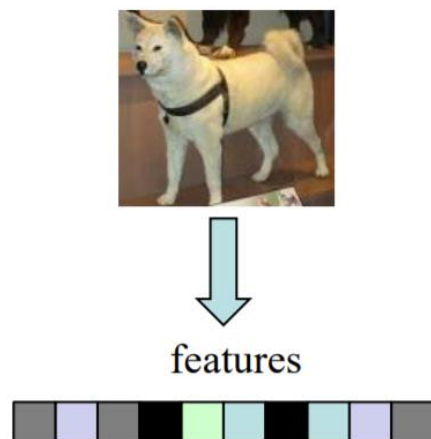
# Introduction (con't)

The key of fine-grained image analysis





# Why hashing



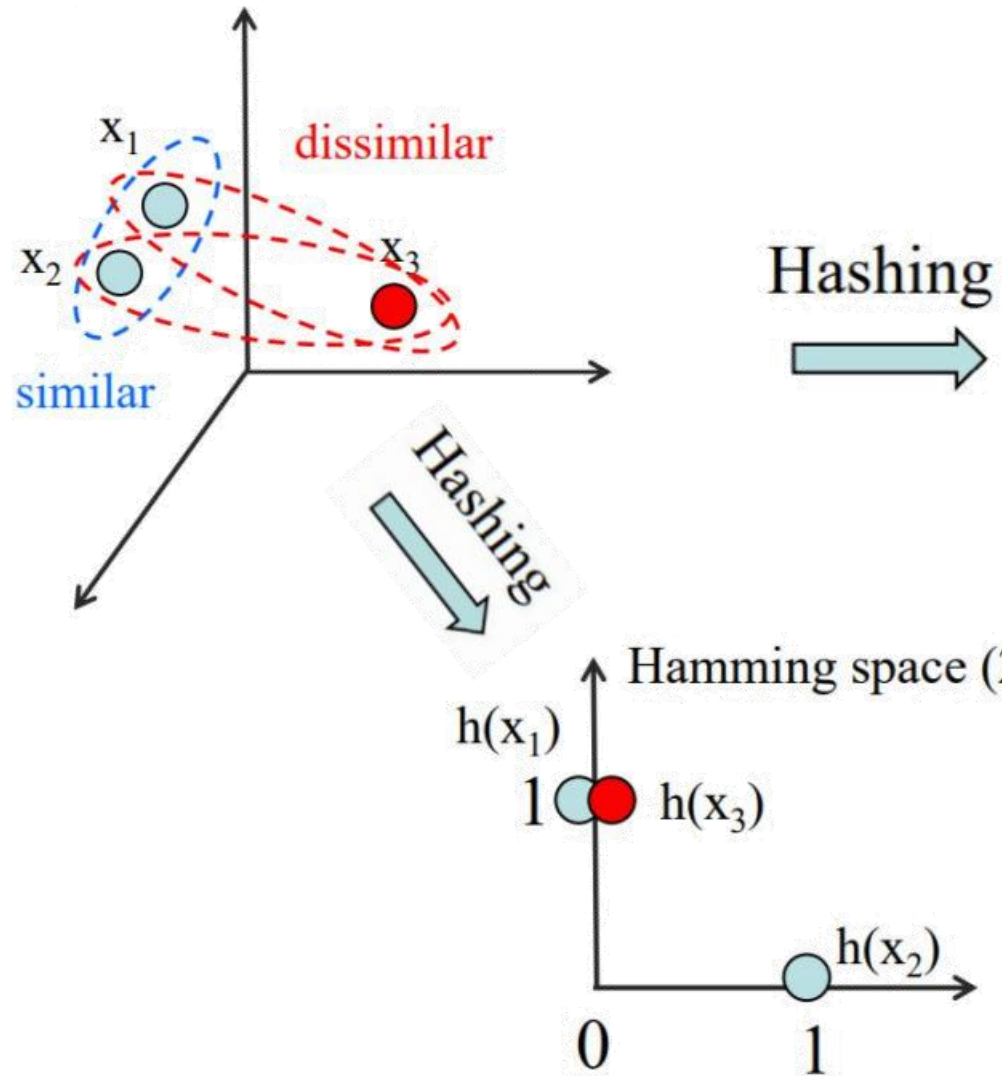
Nearest neighbors search

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.



# Hashing



Hamming space (2-dim) MindSpore

Similarities are well preserved

Similarities are not preserved

# Hashing (con't)

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]

Learn compact  
codes by using  
label information.

Semi-supervised  
Hashing

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

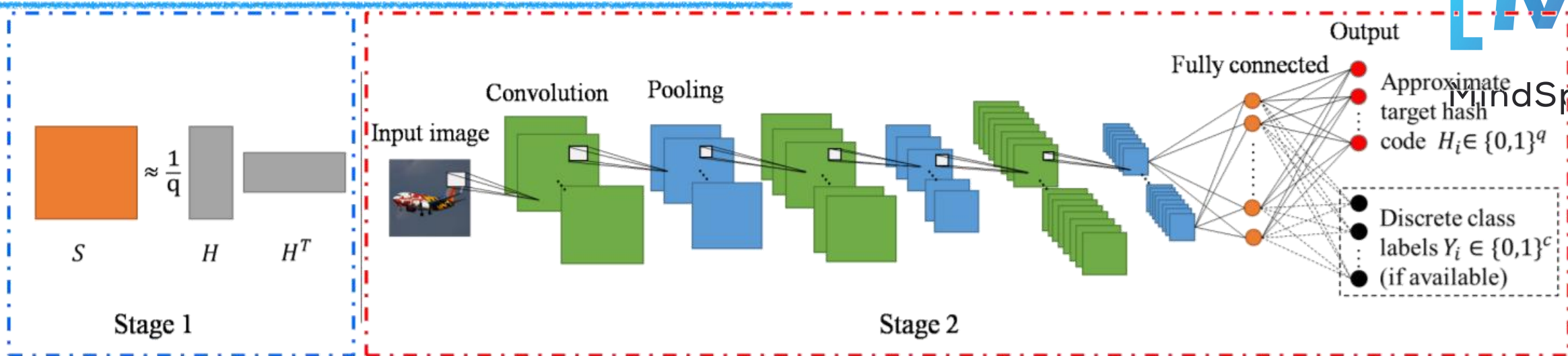
Supervised  
Hashing

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

# Deep Hashing

[M]<sup>s</sup>

MindSpore



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 **ith row in  $H$**  is the approximate hash code for the **ith training image**

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

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



# Fine-grained hashing

## Fine-grained image retrieval

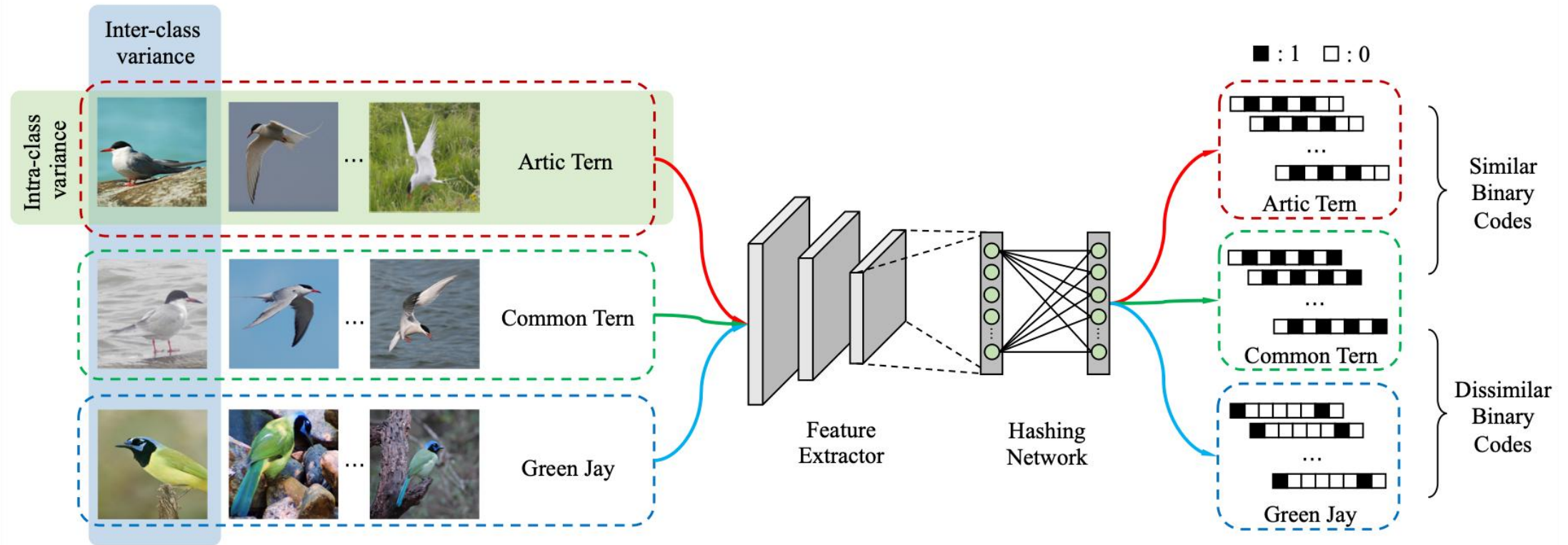


Illustrations of fine-grained image retrieval. Two examples (“Mallard” and “Rolls-Royce Phantom Sedan 2012”) from the *CUB200-2011* and *Cars* datasets, respectively.

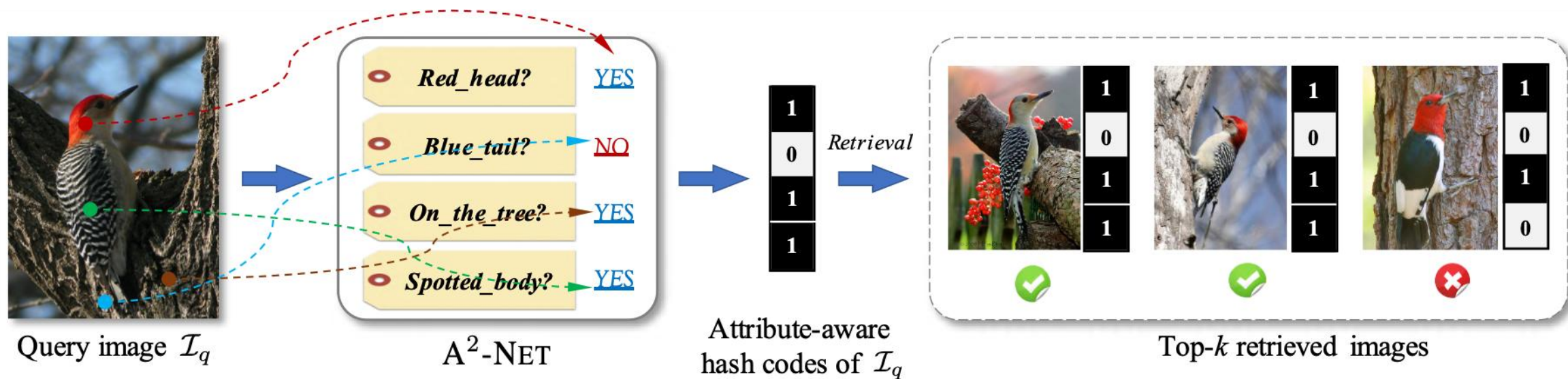


# Fine-grained hashing (con't)

Large-scale fine-grained retrieval requires efficiency ...

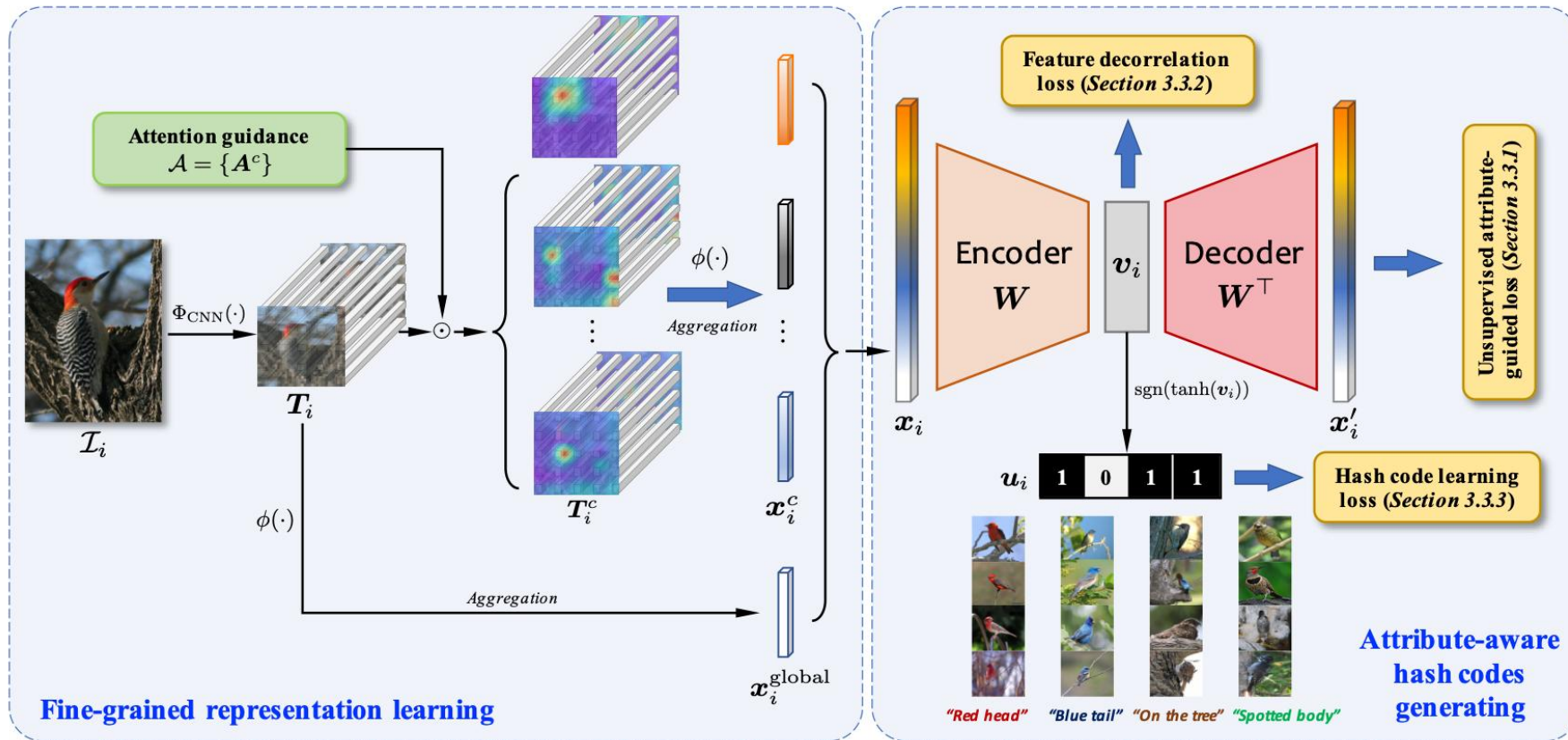


Attribute-aware correspondences are expected ...



# Our works (con't)

## A<sup>2</sup>-Net (Atribute-Aware hashing Network)



$$\min_{\mathbf{W}, \Theta} \mathcal{L}(\mathcal{I}) = \|\mathbf{X} - \mathbf{W}^\top \mathbf{V}'\|_F^2 + \lambda \|\mathbf{W} \mathbf{X} - \mathbf{V}'\|_F^2 + \alpha \|\mathbf{V}' \mathbf{V}'^\top - n \mathbf{I}\|_F^2$$

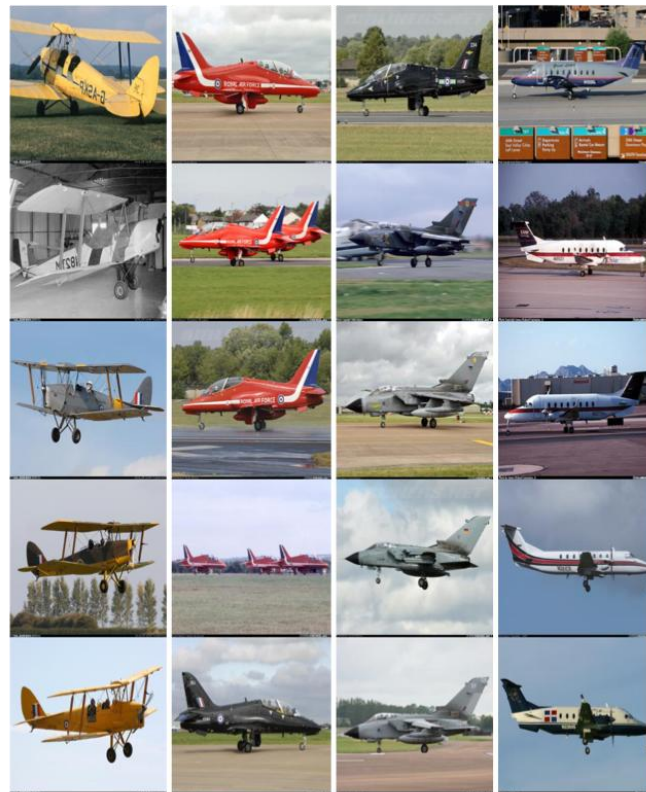
$$+ \beta \sum_{i \in \Omega} \sum_{j \in \Gamma} \left( \tanh(\mathbf{W} \cdot F(\mathcal{I}_i; \Theta))^\top \mathbf{z}_j - k S_{ij} \right)^2.$$

A<sup>2</sup>-NET: Learning Attribute-Aware Hash Codes for Large-Scale Fine-Grained Image Retrieval



# Our works (con't)

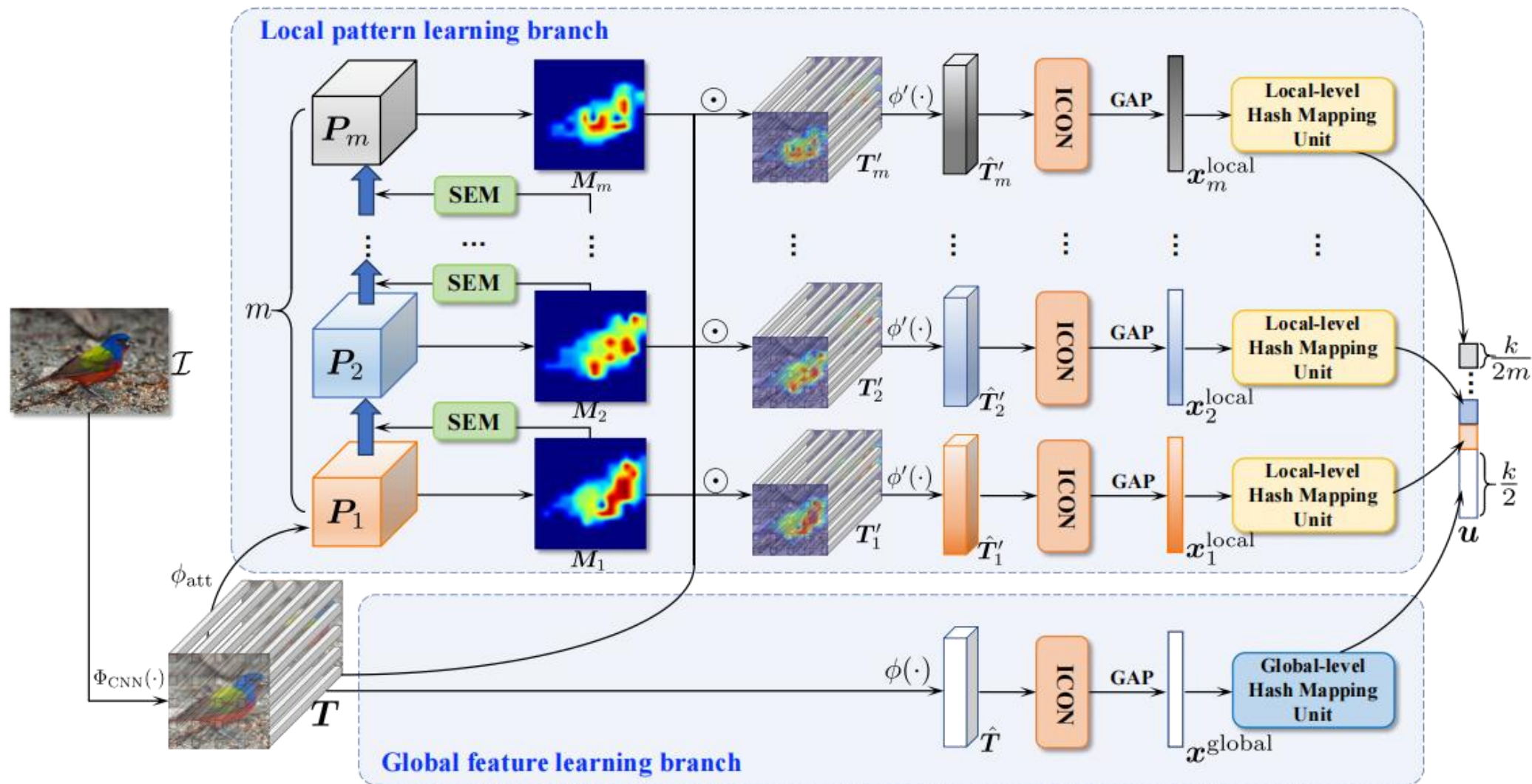
## Quality demonstrations of the learned hash codes



*Attribute-aware!*



# Our works (ECCV 2022)

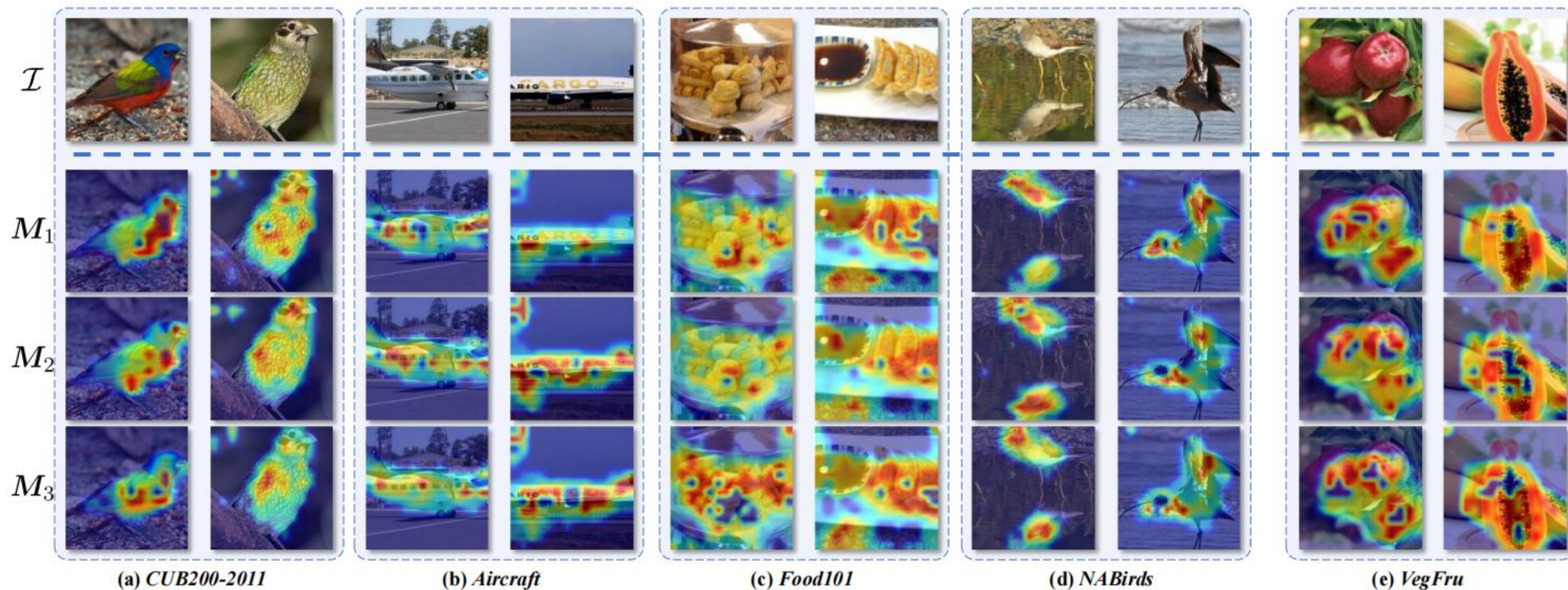




# Our works (con't)

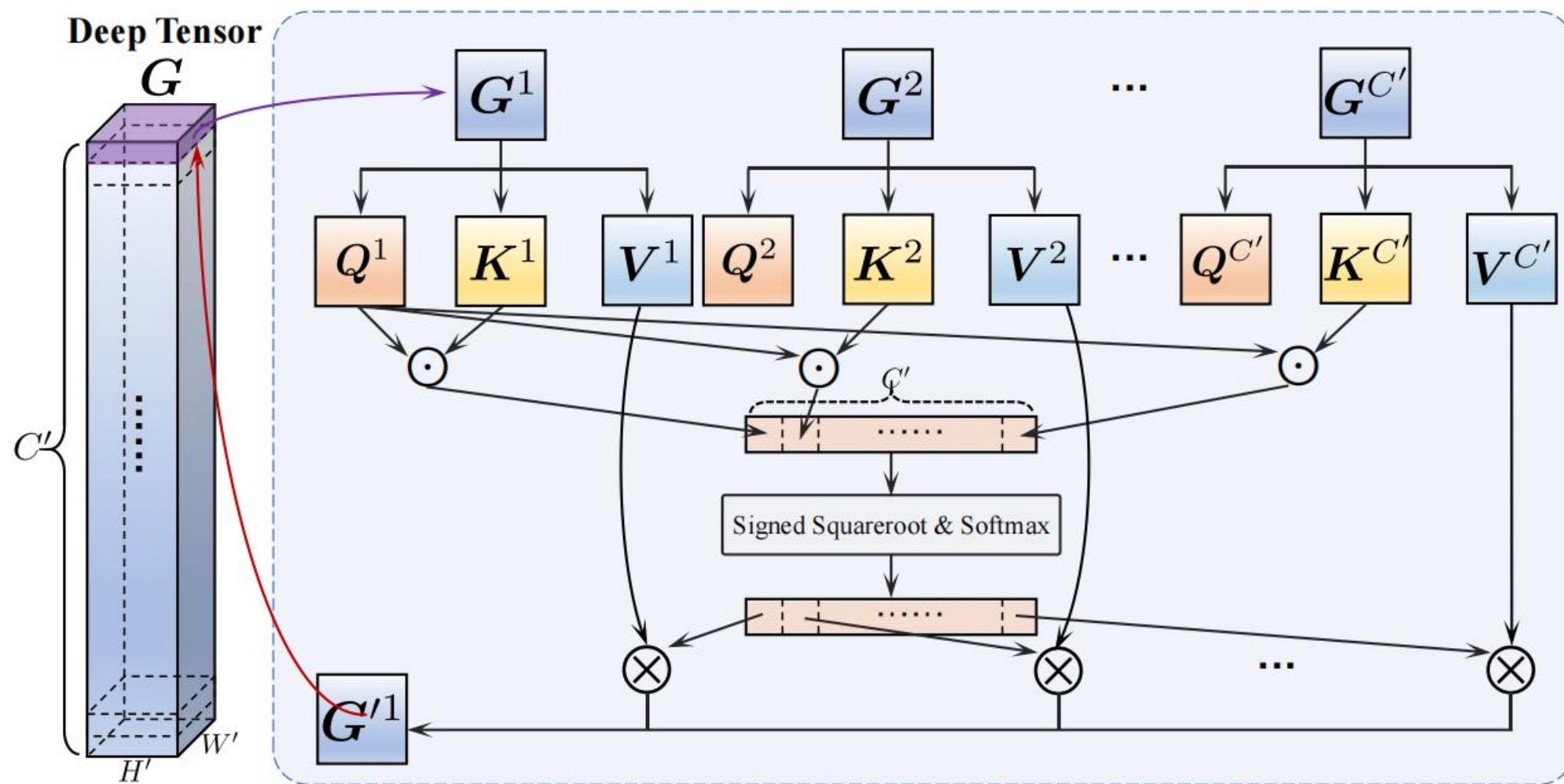


MindSpore



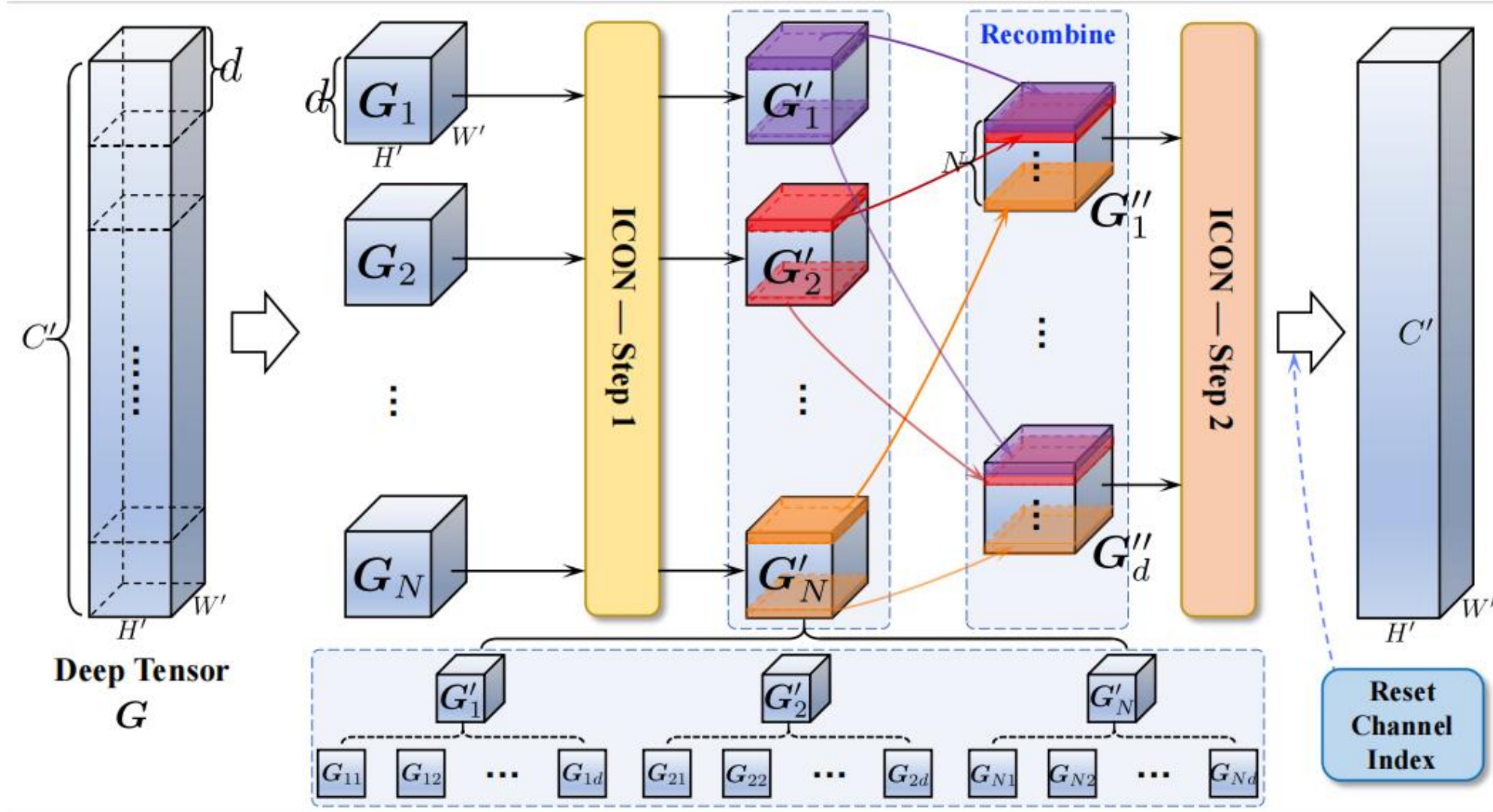
SEMICON: A Learning-to-hash Solution for Large-scale Fine-grained Image Retrieval

# Our works (con't)





# Our works (con't)





# Code with MindSpore



```
class CUB200(object):
    def __init__(self, root, train=True):
        self.root = root
        self.train = train
        image_path = os.path.join(root, 'images')
        id2name = np.genfromtxt(os.path.join(root, 'images.txt'), dtype=str)
        id2train = np.genfromtxt(os.path.join(root, 'train_test_split.txt'), dtype=int)
        self.images = []
        self.labels = []
        target = 1 if train else 0
        for i in range(id2name.shape[0]):
            if id2train[i, 1] != target:
                continue
            label = int(id2name[i, 1][:3]) - 1
            self.images.append(os.path.join(image_path, id2name[i, 1]))
            self.labels.append(label)

    @staticmethod
    def smallest_max_size(x, max_size):
        def py3round(number):
            if abs(round(number) - number) == 0.5:
                return int(2.0 * round(number / 2.0))
            return int(round(number))

        h, w = x.shape[:2]
        scale = max_size / float(min(h, w))
        h, w = tuple(py3round(dim * scale) for dim in (h, w))
        x = cv2.resize(x, (h, w), interpolation=cv2.INTER_LINEAR)
        return x
```

MindSpore 1.3

`mindspore.dataset.vision.c_transforms.AutoContrast`

`mindspore.dataset.vision.c_transforms.BoundingBoxAugment`

`mindspore.dataset.vision.c_transforms.CenterCrop`

`mindspore.dataset.vision.c_transforms.ConvertColor`

`mindspore.dataset.vision.c_transforms.Crop`

`mindspore.dataset.vision.c_transforms.CutMixBatch`

`mindspore.dataset.vision.c_transforms.CutOut`

`mindspore.dataset.vision.c_transforms.Decode`

`mindspore.dataset.vision.c_transforms.Equalize`

`mindspore.dataset.vision.c_transforms.GaussianBlur`

dSpore

MindSpore 1.7

以ECCV 2022的工作为例，采用MindSpore搭建模型

---

# Code with MindSpore (con't)



MindSpore

```
class ChannelTransformer(nn.Cell):
    def __init__(self, dim, num_heads):
        super().__init__()
        self.num_heads = num_heads
        head_dim = dim // num_heads
        self.scale = head_dim ** -0.5
        self.head_dim = head_dim
        self.norm = nn.BatchNorm2d(dim)
        self.relu = nn.ReLU()
        self.softmax = nn.Softmax(axis=-1)
        self.qkv = nn.Conv2d(dim, dim * 3, 1, group=num_heads)
        self.qkv2 = nn.Conv2d(dim, dim * 3, 1, group=head_dim)

    def construct(self, x):
        B, C, H, W = x.shape
        qkv = self.qkv(x).reshape(B, 3, self.num_heads, self.head_dim, H * W).transpose(1,0,2,3,4)
        q, k, v = qkv[0], qkv[1], qkv[2]
        attn = ms.numpy.matmul(q, k.transpose(0,1,3,2)) * self.scale
        attn = ms.numpy.sign(attn) * ms.numpy.sqrt(ms.numpy.abs(attn) + 1e-5)
        attn = self.softmax(attn)
        x = (ms.numpy.matmul(attn, v).reshape(B, C, H, W) + x).reshape(B, self.num_heads, self.head_dim, H * W)
        y = self.norm(x)
        x = self.relu(y)
        qkv2 = self.qkv2(x).reshape(B, 3, self.head_dim, self.num_heads, H * W).transpose(1,0,2,3,4)
        q, k, v = qkv2[0], qkv2[1], qkv2[2]
        attn = (ms.numpy.matmul(q, k.transpose(0,1,3,2))) * (self.num_heads ** -0.5)
        attn = ms.numpy.sign(attn) * ms.numpy.sqrt(ms.numpy.abs(attn) + 1e-5)
        attn = self.softmax(attn)
        x = ms.numpy.matmul(attn, v).reshape(B, self.head_dim, self.num_heads, H, W).transpose(0,2,1,3,4)
        return x
```

MindSpore

```
class ChannelTransformer(nn.Module):
    def __init__(self, dim, num_heads):
        super().__init__()
        self.num_heads = num_heads
        head_dim = dim // num_heads
        self.scale = head_dim ** -0.5
        self.head_dim = head_dim
        self.norm = nn.BatchNorm2d(dim)
        self.relu = nn.ReLU(inplace=True)
        self.qkv = nn.Conv2d(dim, dim * 3, 1, groups=num_heads)
        self.qkv2 = nn.Conv2d(dim, dim * 3, 1, groups=head_dim)

    def forward(self, x):
        B, C, H, W = x.shape
        qkv = self.qkv(x).reshape(B, 3, self.num_heads, self.head_dim, H * W).transpose(0, 1)
        q, k, v = qkv[0], qkv[1], qkv[2]
        attn = (q @ k.transpose(-2, -1)) * self.scale
        attn = torch.sign(attn) * torch.sqrt(torch.abs(attn) + 1e-5)
        attn = attn.softmax(dim=-1)
        x = ((attn @ v).reshape(B, C, H, W) + x).reshape(B, self.num_heads, self.head_dim, H, W).transpose(1, 2)
        y = self.norm(x)
        x = self.relu(y)
        qkv2 = self.qkv2(x).reshape(B, 3, self.head_dim, self.num_heads, H * W).transpose(0, 1)
        q, k, v = qkv2[0], qkv2[1], qkv2[2]
        attn = (q @ k.transpose(-2, -1)) * (self.num_heads ** -0.5)
        attn = torch.sign(attn) * torch.sqrt(torch.abs(attn) + 1e-5)
        attn = attn.softmax(dim=-1)
        x = (attn @ v).reshape(B, self.head_dim, self.num_heads, H, W).transpose(1, 2).reshape(B, C, H, W)
        return x
```

PyTorch

THANK YOU