# Learning to Compare: Relation Network for Few-Shot Learning

Flood Sung    Yongxin Yang[1]    Li Zhang[1,2]    Tao Xiang[1]    Philip H.S. Torr[2]    Timothy M. Hospedales[3]
[1]Queen Mary University of London    [2]University of Oxford    [3]The University of Edinburgh
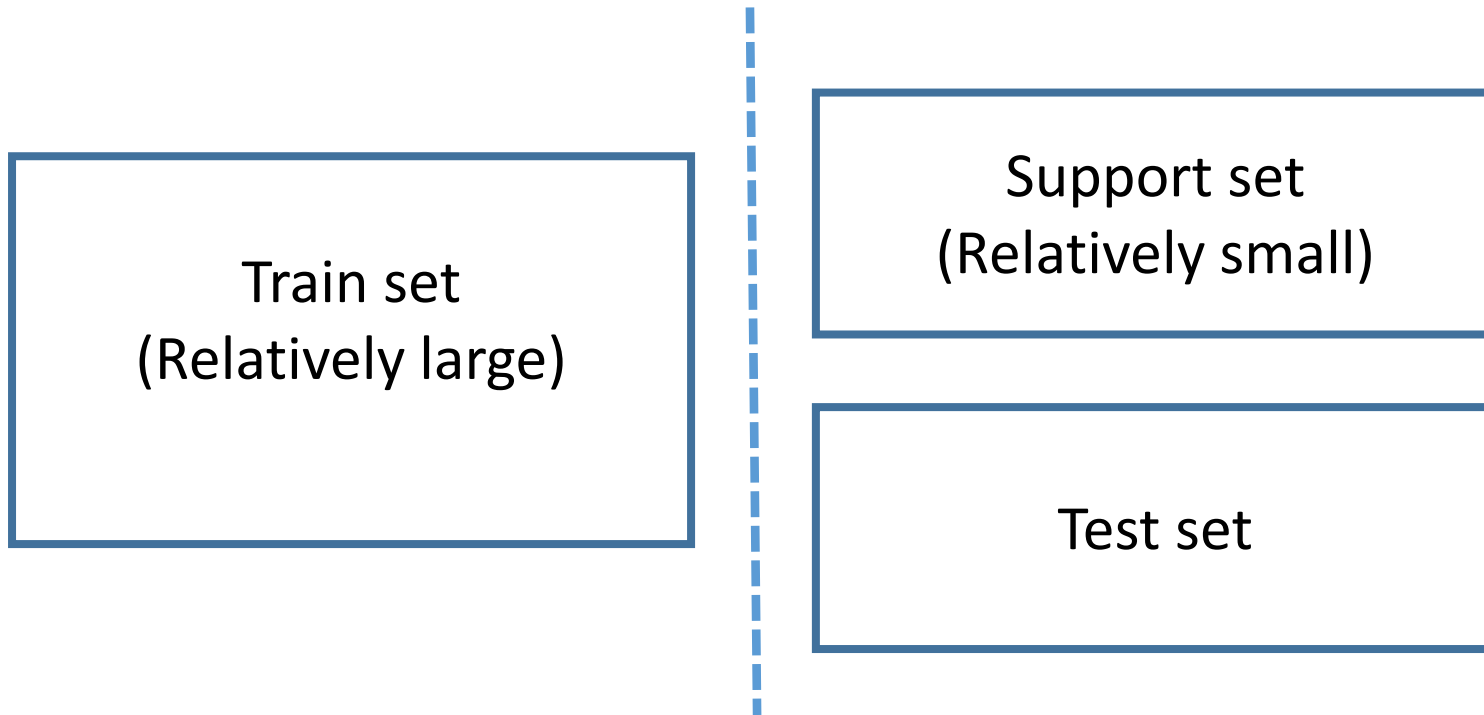floodsung@gmail.com    {yongxin.yang, david.lizhang, t.xiang}@qmul.ac.uk
philip.torr@eng.ox.ac.uk    t.hospedales@ed.ac.uk
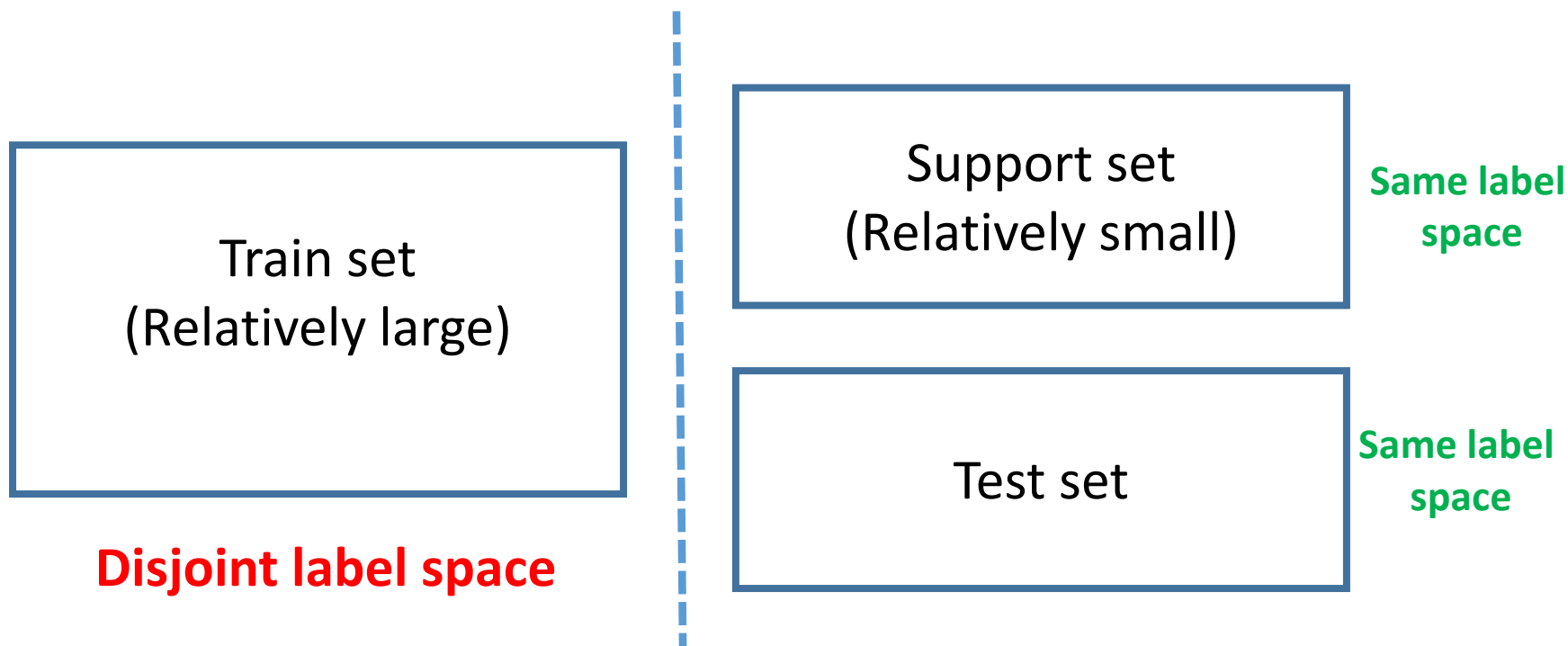
Few-shot learning
Meta-learning

# Few-shot learning

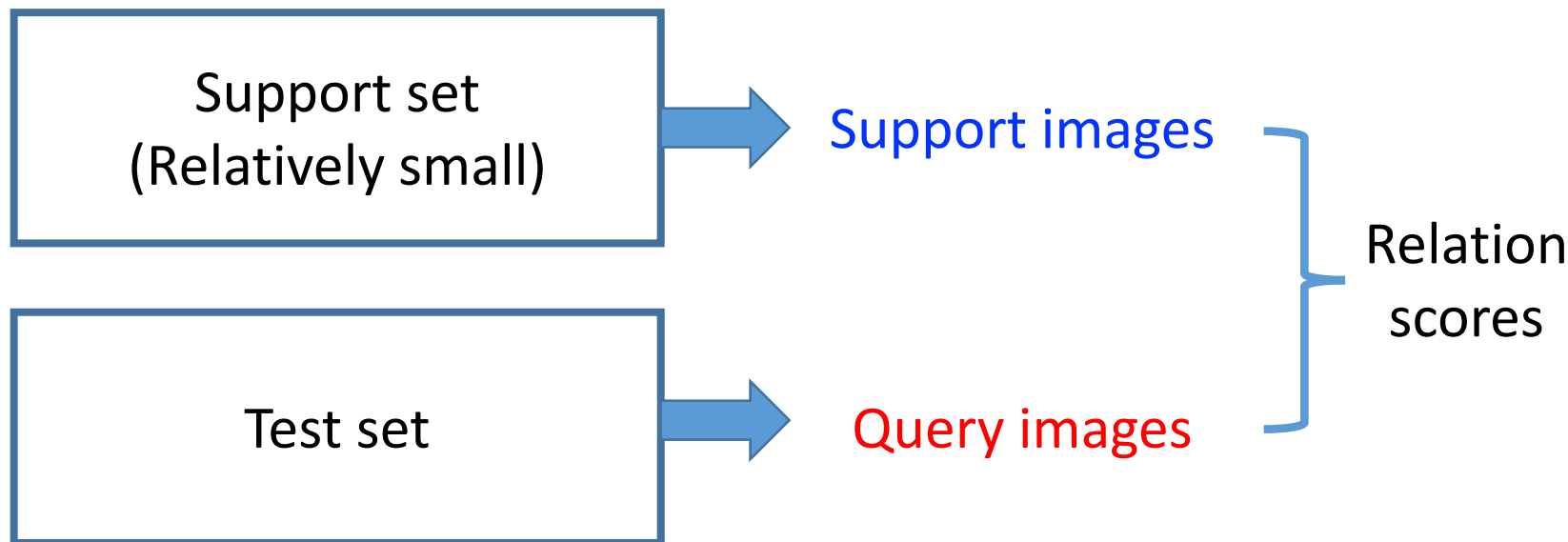Must train a classifier to recognize new classes given few examples from each.

Train set
(Relatively large)

Support set
(Relatively small)

Test set

# Few-shot learning

Must train a classifier to recognize new classes given few examples from each.

| Train set (Relatively large) | | Support set (Relatively small) | Same label space |
|---|---|---|---|
| | | Test set | Same label space |

**Disjoint label space**

# Relation network--classifying by comparison

Relation network (RN) compares <span style="color:blue">support images</span> and <span style="color:red">query images</span> and makes classification according to the returned "**relation scores**"
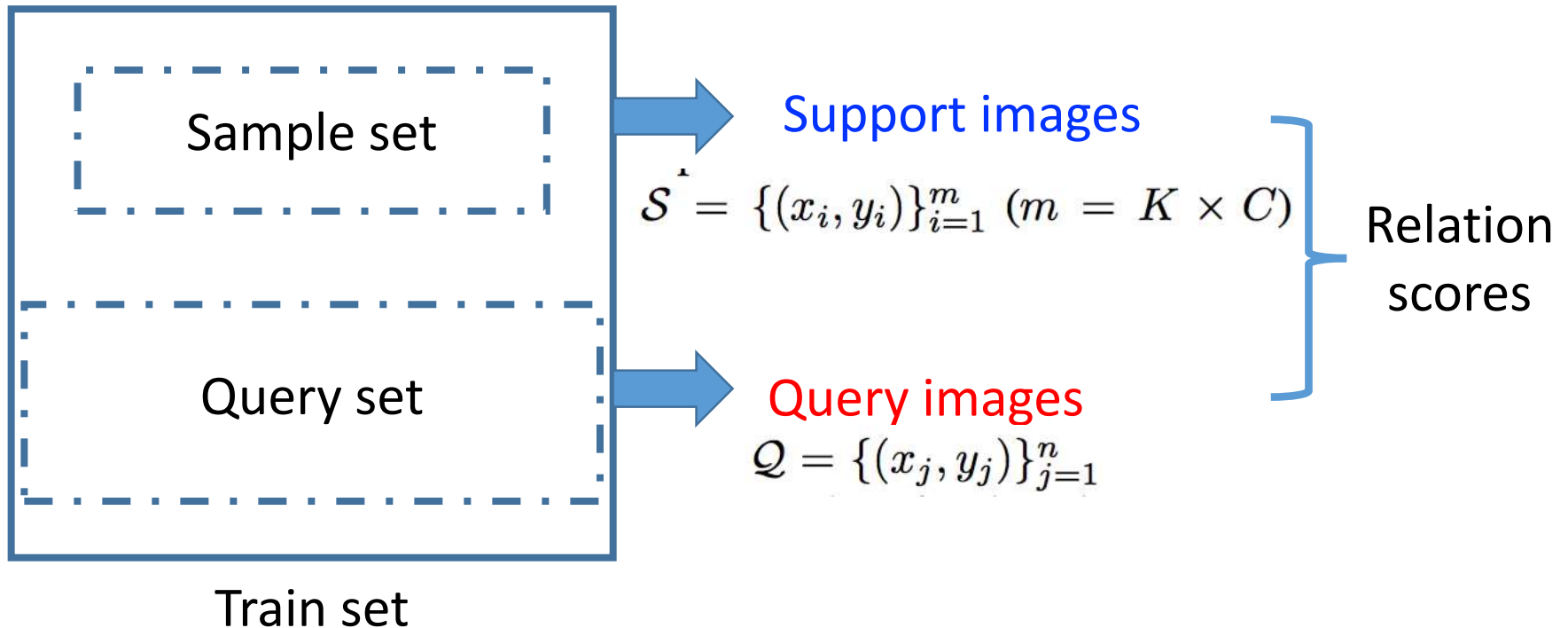


Discussion: Converting classification task into retrieval task?

# Relation network--classifying by comparison

Relation network (RN) learn to compare with meta-learning:
to mimic the comparison procedure on the training set and learn the model.

In each episode

Sample set → Support images

$$\mathcal{S} = \{(x_i, y_i)\}_{i=1}^{m} \ (m = K \times C)$$

Query set → Query images

$$\mathcal{Q} = \{(x_j, y_j)\}_{j=1}^{n}$$

Relation scores

Train set

# Relation network--classifying by comparison

Relation network (RN) learn to compare with meta-learning:
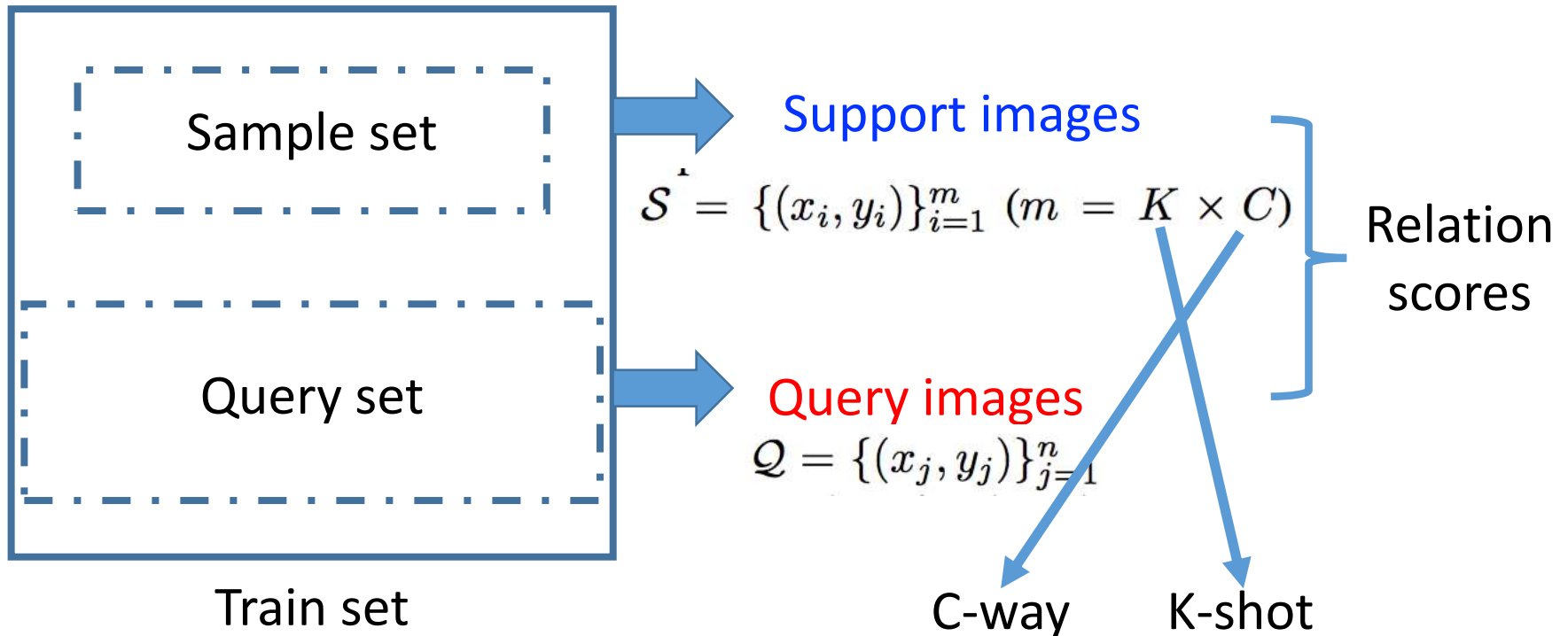to mimic the comparison procedure on the training set and learn
the model.

In each episode

Sample set → Support images

$$\mathcal{S} = \{(x_i, y_i)\}_{i=1}^m \ (m = K \times C)$$

Query set → Query images

$$\mathcal{Q} = \{(x_j, y_j)\}_{j=1}^n$$

Train set

Relation scores
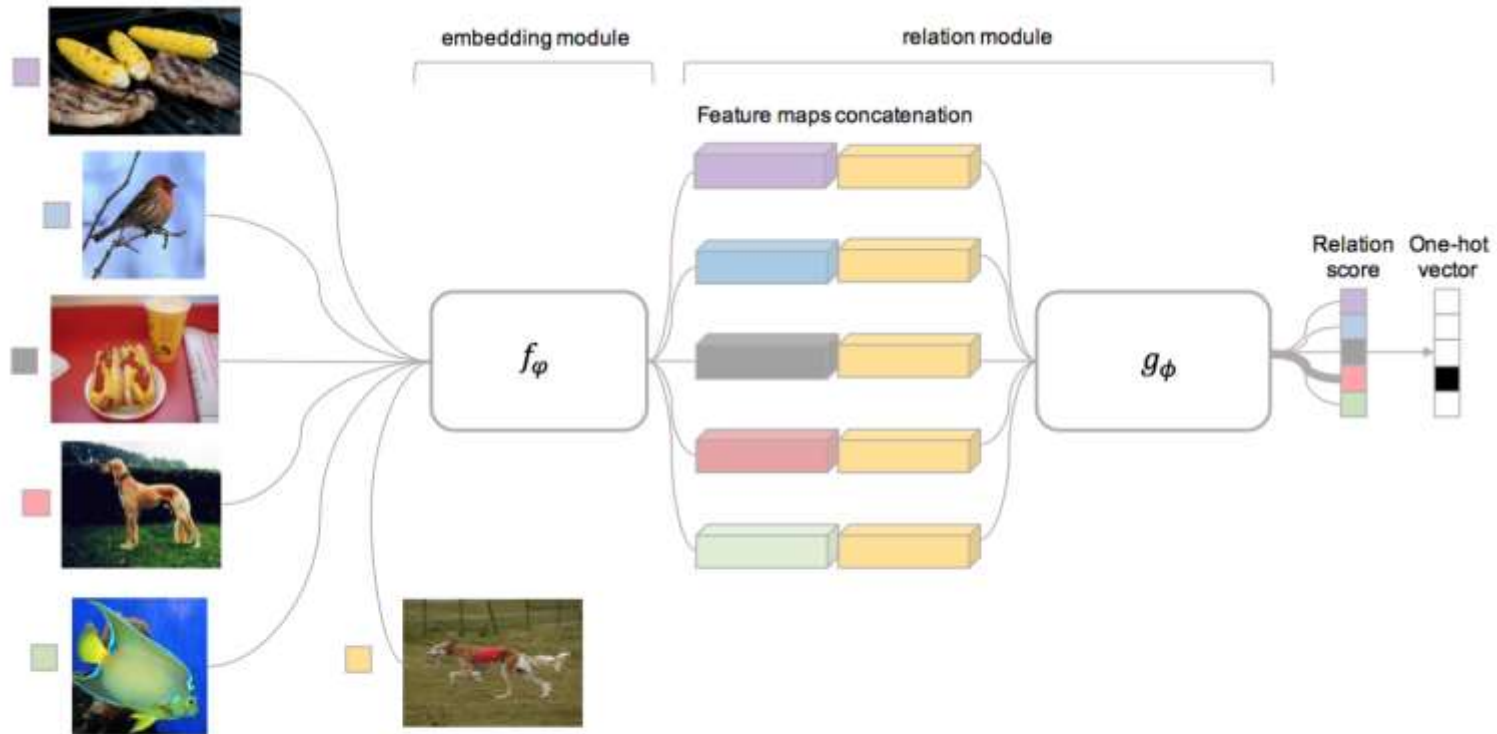
C-way    K-shot

# Relation network--structure



Figure 1: Relation Network architecture with a 5-way 1-shot 1-query example.

$$r_{i,j} = g_\phi(\mathcal{C}(f_\varphi(x_i), f_\varphi(x_j))), \quad i = 1, 2, \ldots, C$$

# Relation network--structure



Figure 1: Relation Network architecture with a 5-way 1-shot 1-query example.

relation module

Embedding module

$$r_{i,j} = g_\phi(\mathcal{C}(f_\varphi(x_i), f_\varphi(x_j))), \quad i = 1, 2, \ldots, C$$

$r_{i,j}$ is bounded between (0,1) by sigmoid function
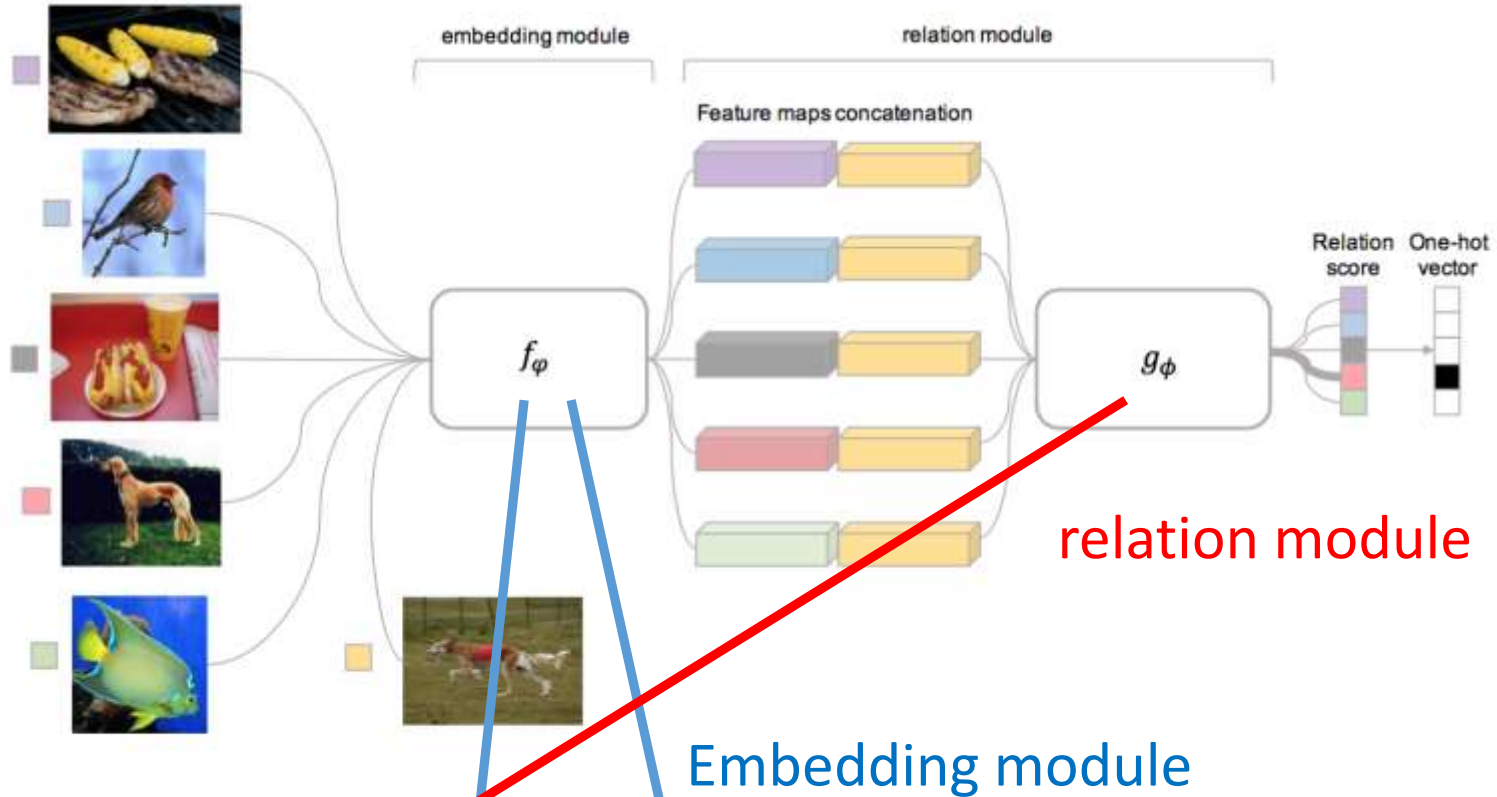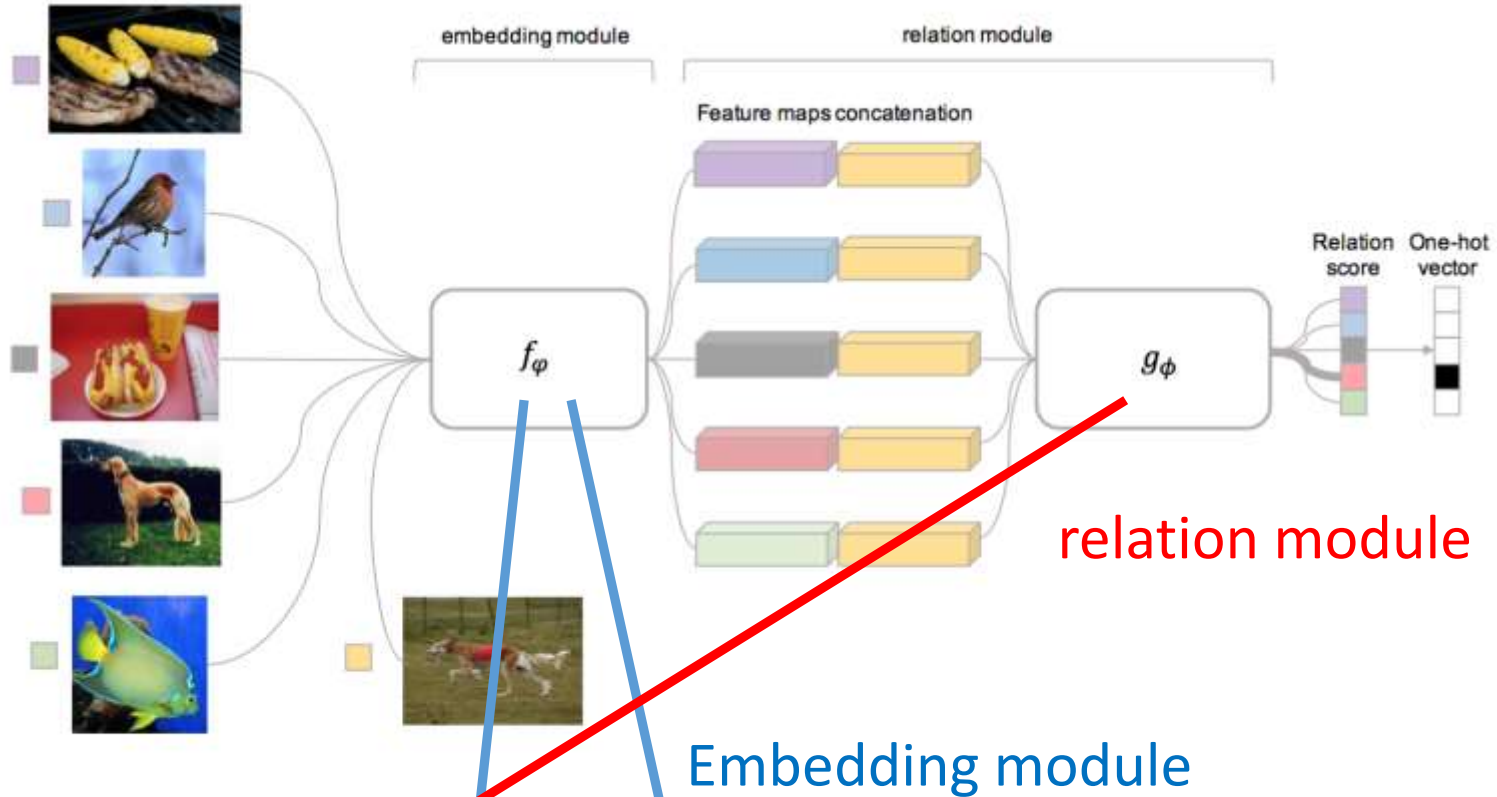
# Relation network--structure



Figure 1: Relation Network architecture with a 5-way 1-shot 1-query example.

$$r_{i,j} = g_\phi(\mathcal{C}(f_\varphi(x_i), f_\varphi(x_j))), \quad i = 1, 2, \ldots, C$$

Optimization target: $\varphi, \phi \leftarrow \underset{\varphi,\phi}{\mathrm{argmin}} \sum_{i=1}^{m} \sum_{j=1}^{n} (r_{i,j} - \mathbf{1}(y_i == y_j))^2$

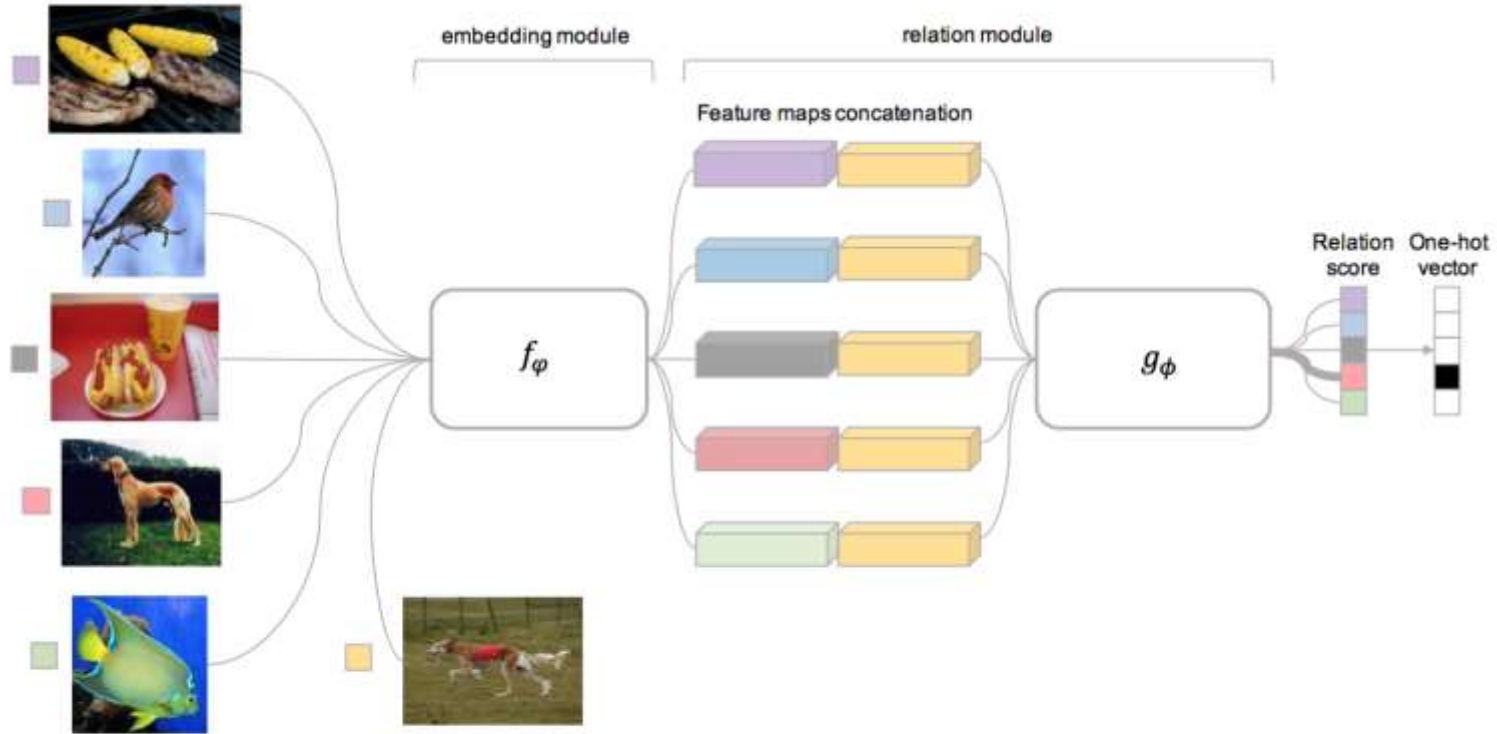# Relation network--structure



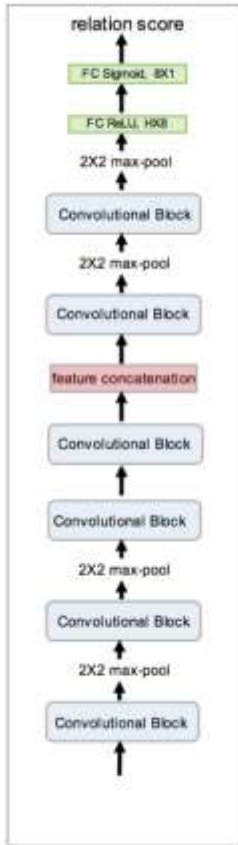Figure 1: Relation Network architecture with a 5-way 1-shot 1-query example.

A detail for K-shot: K embedded features are pooled by pixel-wise sum operation
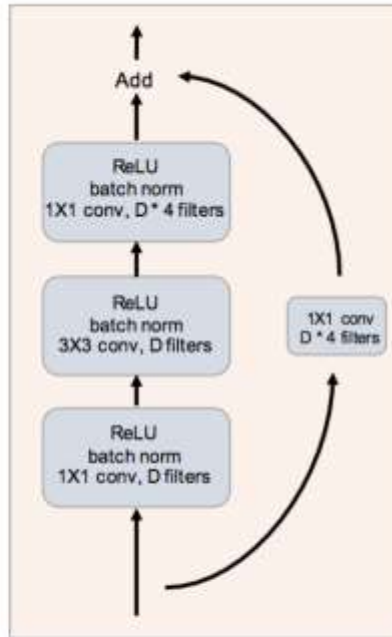
# Relation network--structure
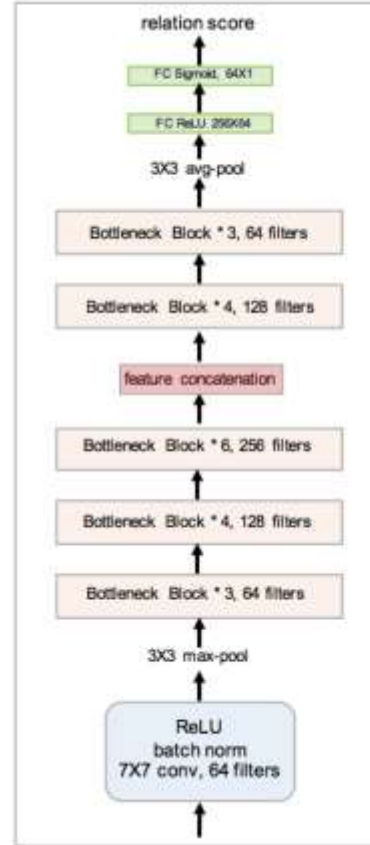


(a) Convolutional Block

ReLU
batch norm
3X3 conv, 64 filters
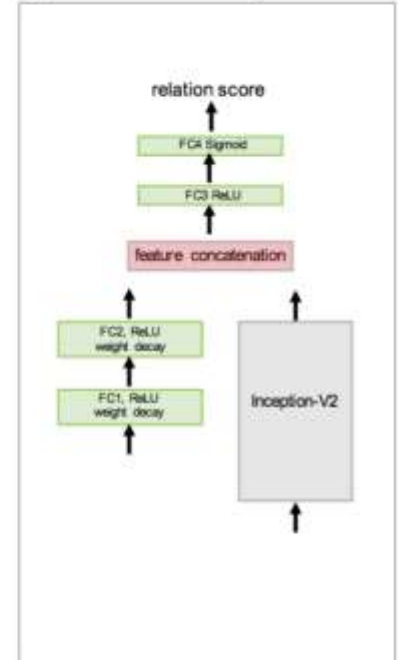
(b) Naive RN for few-shot learning

relation score

FC Sigmoid, 8X1

FC ReLU, HX8

2X2 max-pool

Convolutional Block

2X2 max-pool

Convolutional Block

feature concatenation

Convolutional Block

Convolutional Block

2X2 max-pool

Convolutional Block

2X2 max-pool

Convolutional Block

(c) Bottleneck Block

Add

ReLU
batch norm
1X1 conv, D * 4 filters

ReLU
batch norm
3X3 conv, D filters

1X1 conv
D * 4 filters

ReLU
batch norm
1X1 conv, D filters

(d) Deeper RN for few-shot learning

relation score

FC Sigmoid, 64X1

FC ReLU, 256X64

3X3 avg-pool

Bottleneck Block * 3, 64 filters

Bottleneck Block * 4, 128 filters

feature concatenation

Bottleneck Block * 6, 256 filters

Bottleneck Block * 4, 128 filters

Bottleneck Block * 3, 64 filters

3X3 max-pool

ReLU
batch norm
7X7 conv, 64 filters

(e) RN for zero-shot learning

relation score

FC4 Sigmoid

FC3 ReLU

feature concatenation

FC2, ReLU
weight decay

FC1, ReLU
weight decay

Inception-V2

Detailed structure

Extension to 0-shot learning: different embedding module for sample and query images

Semantic vector        images

# Experiments

| Model | Fine Tune | 5-way Acc. | | 20-way Acc. | |
|---|---|---|---|---|---|
| | | 1-shot | 5-shot | 1-shot | 5-shot |
| MANN [31] | N | 82.8% | 94.9% | - | - |
| CONVOLUTIONAL SIAMESE NETS [18] | N | 96.7% | 98.4% | 88.0% | 96.5% |
| CONVOLUTIONAL SIAMESE NETS [18] | Y | 97.3% | 98.4% | 88.1% | 97.0% |
| MATCHING NETS [38] | N | 98.1% | 98.9% | 93.8% | 98.5% |
| MATCHING NETS [38] | Y | 97.9% | 98.7% | 93.5% | 98.7% |
| SIAMESE NETS WITH MEMORY [16] | N | 98.4% | 99.6% | 95.0% | 98.6% |
| NEURAL STATISTICIAN [8] | N | 98.1% | 99.5% | 93.2% | 98.1% |
| META NETS [26] | N | 99.0% | - | 97.0% | - |
| PROTOTYPICAL NETS [35] | N | 98.8% | 99.7% | 96.0% | 98.9% |
| MAML [10] | Y | $98.7 \pm 0.4\%$ | $99.9 \pm 0.1\%$ | $95.8 \pm 0.3\%$ | $98.9 \pm 0.2\%$ |
| RELATION NET | N | $99.6 \pm 0.2\%$ | $99.8 \pm 0.1\%$ | $97.6 \pm 0.2\%$ | $99.1 \pm 0.1\%$ |

Table 1: Omniglot few-shot classification. Results are accuracies averaged over 1000 test episodes and with 95% confidence intervals where reported. The best-performing method is highlighted, along with others whose confidence intervals overlap. '-': not reported.
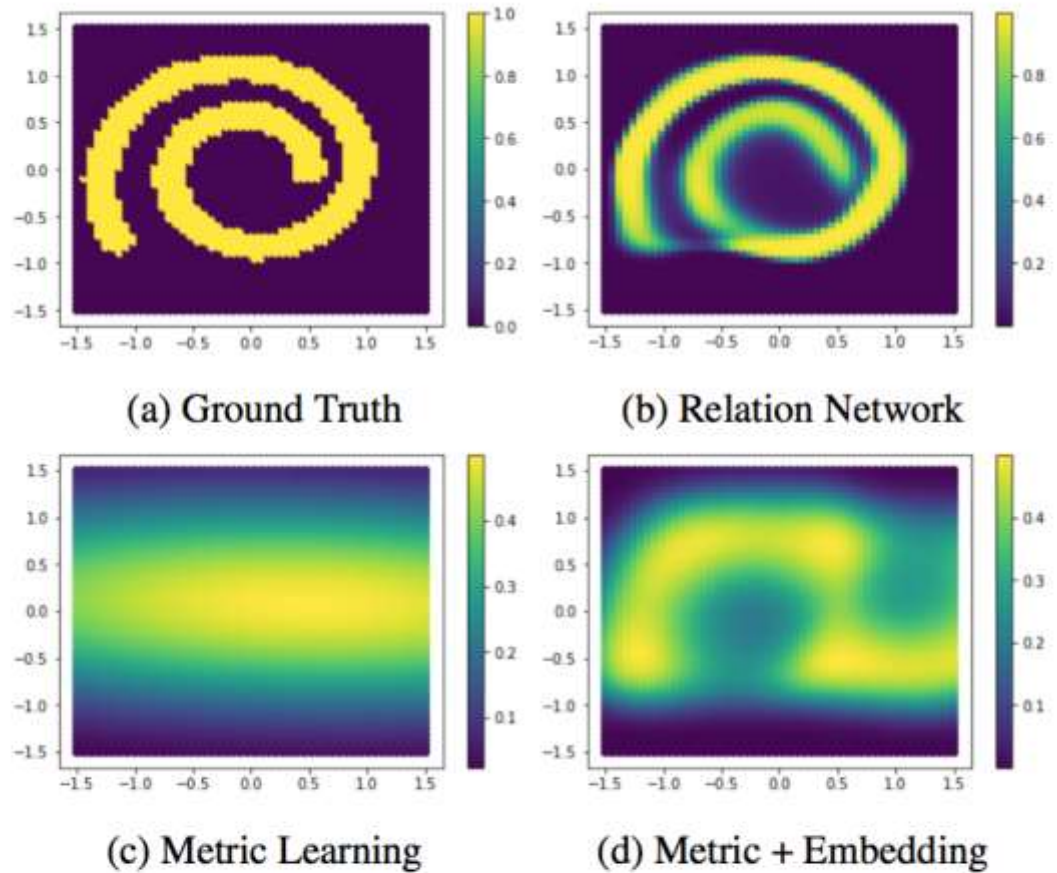
# Why does RN work?

Both deep feature embedding and deep <span style="color:red">distance metric</span> are learnable.

The concatenation operation is in relatively <span style="color:red">bottom layer</span>

When training a Siamese network or a triplet network, we apply metric constraint on a specified feature and then use the Euclidean distance (or other fixed metric) for metric for inference.
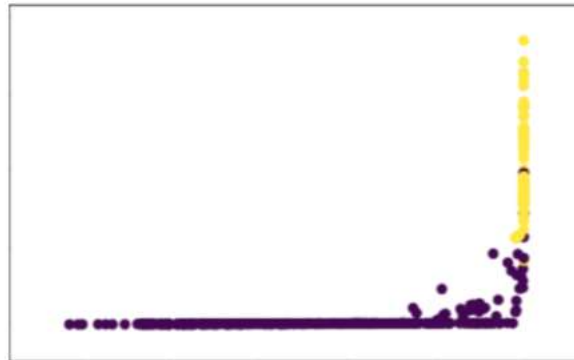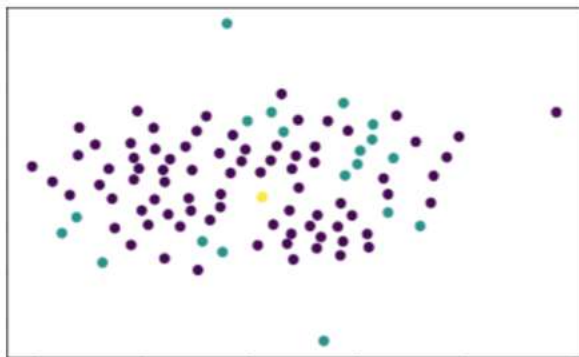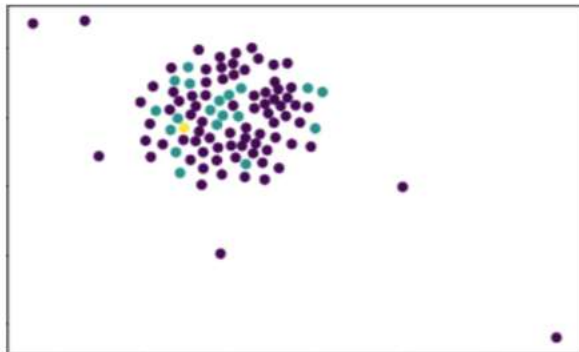
# Why does RN work?



2D data space

(a) Ground Truth

(b) Relation Network

(c) Metric Learning

(d) Metric + Embedding

Figure 3: An example relation learnable by Relation Network and not by non-linear embedding + metric learning.

# Why does RN work?



The feature embeddings are difficult to separate.

The relation module pair representations are linearly separable

Figure 4: Example Omniglot few-shot problem visualisations. Left: Matched (cyan) and mismatched (magenta) sample embeddings for a given query (yellow) are not straightforward to differentiate. Right: Matched (yellow) and mismatched (magenta) relation module pair representations are linearly separable.

# Why does RN work? My guess

1) Feature concatenation operation in very early stage (bottom layers)

2) K sample images to mimic the K-shot

3) Converting classification to "comparison", which is a semi-parameter model approach.