

Learning to Compare: Relation Network for Few-Shot Learning

Flood Sung¹ Yongxin Yang¹ Li Zhang^{1,2} Tao Xiang¹ Philip H.S. Torr² Timothy M. Hospedales³
¹Queen Mary University of London ²University of Oxford ³The 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

Abstract

We present a conceptually simple, flexible, and general framework for few-shot learning, where a classifier must learn to recognise new classes given only few examples from each. Our method, called the Relation Network (RN), is trained end-to-end from scratch. During meta-learning, it learns to learn a deep distance metric to compare a small number of images within episodes, each of which is designed to simulate the few-shot setting. Once trained, a RN is able to classify images of new classes by computing relation scores between query images and the few examples of each new class without further updating the network. Besides providing improved performance on few-shot learning, our framework is easily extended to zero-shot learning. Extensive experiments on four datasets demonstrate that our simple approach provides a unified and effective approach for both of these two tasks.

1. Introduction

Deep learning models have achieved great success in visual recognition tasks [20, 14, 34]. However, these supervised learning models need large amounts of labelled data and many iterations to train their large number of parameters. This severely limits their scalability to new classes due to annotation cost, but more fundamentally limits their applicability to newly emerging (eg. new consumer devices) or rare (eg. rare animals) categories where numerous annotated images may simply never exist. In contrast, humans are very good at recognising objects with very little direct supervision, or none at all *i.e.*, few-shot [21, 9] or zero-shot [22] learning. For example, children have no problem generalising the concept of “zebra” from a single picture in a book, or hearing its description as looking like a stripy horse. Motivated by the failure of conventional deep learning methods to work well on one or few examples per class, and inspired by the few- and zero-shot learning ability of humans, there has been a recent resurgence of interest in machine one/few-shot [8, 38, 31, 16, 18, 10, 26, 35, 28] and

zero-shot [11, 2, 22, 41, 23, 30] learning.

Few-shot learning aims to recognise novel visual categories from very few labelled examples. The availability of only one or very few examples challenges the standard ‘fine-tuning’ practice in deep learning [10]. Data augmentation and regularisation techniques can alleviate overfitting in such a limited-data regime, but they do not solve it. Therefore contemporary approaches to few-shot learning often decompose training into an auxiliary meta learning phase where transferrable knowledge is learned in the form of good initial conditions [10], embeddings [35, 38] or optimisation strategies [28]. The target few-shot learning problem is then learned by fine-tuning [10] with the learned optimisation strategy [28] or computed in a feed-forward pass [35, 38, 3, 31] without updating network weights. Zero-shot learning also suffers from a related challenge. Recognisers are trained by a single example in the form of a class description (c.f., single exemplar image in one-shot), making data insufficiency for gradient-based learning a challenge.

While promising, most existing few-shot learning approaches either require complex inference mechanisms [21, 9], complex recurrent neural network (RNN) architectures [38, 31], or fine-tuning the target problem [10, 28]. Our approach is most related to others that aim to train an effective metric for one-shot learning [38, 35, 18]. Where they focus on the learning of the transferrable embedding and pre-define a fixed metric (e.g., as Euclidean [35]), we further aim to *learn* a transferrable deep metric for comparing the relation between images (few-shot learning), or between images and class descriptions (zero-shot learning). By expressing the inductive bias of a *deeper* solution (multiple non-linear learned stages at both embedding and relation modules), we make it easier to learn a generalisable solution to the problem.

Specifically, we propose a two-branch Relation Network (RN) that performs few-shot recognition by learning to compare *query* images against few-shot labeled *sample* images. First an *embedding module* generates representations of the query and training images. Then these embeddings are compared by a *relation module* that determines if they

are from matching categories or not. Defining an episode-based strategy inspired by [38, 35], the embedding and relation modules are meta-learned end-to-end to support few-shot learning. This can be seen as extending the strategy of [38, 35] to include a *learnable non-linear* comparator, instead of a fixed linear comparator. Our approach outperforms prior approaches, while being simpler (no RNNs [38, 31, 28]) and faster (no fine-tuning [28, 10]). Our proposed strategy also directly generalises to zero-shot learning. In this case the sample branch embeds a single-shot category description rather than a single exemplar training image, and the relation module learns to compare query image and category description embeddings.

Overall our contribution is to provide a clean framework that elegantly encompasses both few and zero-shot learning. Our evaluation on four benchmarks show that it provides compelling performance across the board while being simpler and faster than the alternatives.

2. Related Work

The study of one or few-shot object recognition has been of interest for some time [9]. Earlier work on few-shot learning tended to involve generative models with complex iterative inference strategies [9, 21]. With the success of discriminative deep learning-based approaches in the data-rich many-shot setting [20, 14, 34], there has been a surge of interest in generalising such deep learning approaches to the few-shot learning setting. Many of these approaches use a meta-learning or learning-to-learn strategy in the sense that they extract some transferrable knowledge from a set of auxiliary tasks (meta-learning, learning-to-learn), which then helps them to learn the target few-shot problem well without suffering from the overfitting that might be expected when applying deep models to sparse data problems.

Learning to Fine-Tune The successful MAML approach [10] aimed to meta-learn an initial condition (set of neural network weights) that is good for fine-tuning on few-shot problems. The strategy here is to search for the weight configuration of a given neural network such that it can be effectively fine-tuned on a sparse data problem within a few gradient-descent update steps. Many distinct target problems are sampled from a multiple task training set; the base neural network model is then fine-tuned to solve each of them, and the success at each target problem after fine-tuning drives updates in the base model – thus driving the production of an easy to fine-tune initial condition. The few-shot optimisation approach [28] goes further in meta-learning not only a good initial condition but an LSTM-based optimizer that is trained to be specifically effective for fine-tuning. However both of these approaches suffer from the need to fine-tune on the target problem. In contrast, our approach solves target problems in an entirely feed-forward

manner with no model updates required, making it more convenient for low-latency or low-power applications.

RNN Memory Based Another category of approaches leverage recurrent neural networks with memories [26, 31]. Here the idea is typically that an RNN iterates over an examples of given problem and accumulates the knowledge required to solve that problem in its hidden activations, or external memory. New examples can be classified, for example by comparing them to historic information stored in the memory. So ‘learning’ a single target problem can occur in unrolling the RNN, while learning-to-learn means training the weights of the RNN by learning many distinct problems. While appealing, these architectures face issues in ensuring that they reliably store all the, potentially long term, historical information of relevance without forgetting. In our approach we avoid the complexity of recurrent networks, and the issues involved in ensuring the adequacy of their memory. Instead our learning-to-learn approach is defined entirely with simple and fast feed forward CNNs.

Embedding and Metric Learning Approaches The prior approaches entail some complexity when learning the target few-shot problem. Another category of approach aims to learn a set of projection functions that take query and sample images from the target problem and classify them in a feed forward manner [38, 35, 3]. One approach is to parameterise the weights of a feed-forward classifier in terms of the sample set [3]. The meta-learning here is to train the auxiliary parameterisation net that learns how to parameterise a given feed-forward classification problem in terms of a few-shot sample set. Metric-learning based approaches aim to learn a set of projection functions such that when represented in this embedding, images are easy to recognise using simple nearest neighbour or linear classifiers [38, 35, 18]. In this case the meta-learned transferrable knowledge are the projection functions and the target problem is a simple feed-forward computation.

The most related methodologies to ours are the prototypical networks of [35] and the siamese networks of [18]. These approaches focus on learning embeddings that transform the data such that it can be recognised with a *fixed* nearest-neighbour [35] or linear [18, 35] classifier. In contrast, our framework further defines a relation classifier CNN, in the style of [32] (While [32] focuses on reasoning about relation between two objects in a same image which is to address a different problem.). Compared to [18, 35], this can be seen as providing a learnable rather than fixed metric, or non-linear rather than linear classifier. Compared to [18] we benefit from an episodic training strategy with an end-to-end manner from scratch, and compared to [31] we avoid the complexity of set-to-set RNN embedding of the sample-set, and simply rely on pooling [32].

Zero-Shot Learning Our approach is designed for few-

shot learning, but elegantly spans the space into zero-shot learning (ZSL) by modifying the sample branch to input a single category description rather than single training image. When applied to ZSL our architecture is related to methods that learn to align images and category embeddings and perform recognition by predicting if an image and category embedding pair match [11, 2, 40, 42]. Similarly to the case with the prior metric-based few-shot approaches, most of these apply a fixed manually defined similarity metric or linear classifier after combining the image and category embedding. In contrast, we again benefit from a deeper end-to-end architecture including a learned non-linear metric in the form of our learned convolutional relation network; as well as from an episode-based training strategy.

3. Methodology

3.1. Problem definition

We consider the task of few-shot classifier learning. Formally, we have three datasets: a training set, a support set, and a testing set. The support set and testing set share the same label space, but the training set has its own label space that is disjoint with support/testing set. If the support set contains K labelled examples for each of C unique classes, the target few-shot problem is called C -way K -shot.

With the support set only, we can in principle train a classifier to assign a class label \hat{y} to each sample \hat{x} in the test set. However, due to the lack of labelled samples in the support set, the performance of such a classifier is usually not satisfactory. Therefore we aim to perform meta-learning on the training set, in order to extract transferrable knowledge that will allow us to perform better few-shot learning on the support set and thus classify the test set more successfully.

An effective way to exploit the training set is to mimic the few-shot learning setting via *episode* based training, as proposed in [38]. In each training iteration, an episode is formed by randomly selecting C classes from the training set with K labelled samples from each of the C classes to act as the *sample* set $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^m$ ($m = K \times C$), as well as a fraction of the remainder of those C classes' samples to serve as the *query* set $\mathcal{Q} = \{(x_j, y_j)\}_{j=1}^n$. That this sample/query set split is designed to simulate the support/test set that will be encountered at test time. A model trained from sample/query set can be further fine-tuned using the support set, if desired. In this work we adopt such an episode-based training strategy. In our few-shot experiments (see Section 4.1) we consider one-shot ($K = 1$, Figure 1) and five-shot ($K = 5$) settings. We also address the $K = 0$ zero-shot learning case as explained in Section 3.3.

3.2. Model

One-Shot Our Relation Network (RN) consists of two modules: an *embedding* module f_ϕ and a *relation* module g_ϕ , as illustrated in Figure 1. Samples x_j in the query set \mathcal{Q} , and samples x_i in the sample set \mathcal{S} are fed through the embedding module f_ϕ , which produces feature maps $f_\phi(x_i)$ and $f_\phi(x_j)$. The feature maps $f_\phi(x_i)$ and $f_\phi(x_j)$ are combined with operator $\mathcal{C}(f_\phi(x_i), f_\phi(x_j))$. In this work we assume $\mathcal{C}(\cdot, \cdot)$ to be concatenation of feature maps in depth, although other choices are possible.

The combined feature map of the sample and query are fed into the relation module g_ϕ , which eventually produces a scalar in range of 0 to 1 representing the similarity between x_i and x_j , which is called relation score. Thus, in the C -way one-shot setting, we generate C relation scores $r_{i,j}$ for the relation between one query input x_j and training sample set examples x_i ,

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

K-shot For K -shot where $K > 1$, we element-wise sum over the embedding module outputs of all samples from each training class to form this class' feature map. This pooled class-level feature map is combined with the query image feature map as above. Thus, the number of relation scores for one query is always C in both one-shot or few-shot setting.

Objective function We use mean square error (MSE) loss (Eq. (2)) to train our model, regressing the relation score $r_{i,j}$ to the ground truth: matched pairs have similarity 1 and the mismatched pair have similarity 0.

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

The choice of MSE is somewhat non-standard. Our problem may seem to be a classification problem with a label space $\{0, 1\}$. However conceptually we are predicting similarity scores, which can be considered a regression problem despite that for ground-truth we can only automatically generate $\{0, 1\}$ targets.

3.3. Zero-shot learning

Zero-shot learning is analogous to one-shot learning in that one datum is given to define each class to recognise. However instead of being given a support set with one-shot image for each of C training classes, it contains a semantic class embedding vector v_c for each. Modifying our framework to deal with the zero-shot case is straightforward: as a different modality of semantic vectors is used for the *support* set (e.g. attribute vectors instead of images), we use a

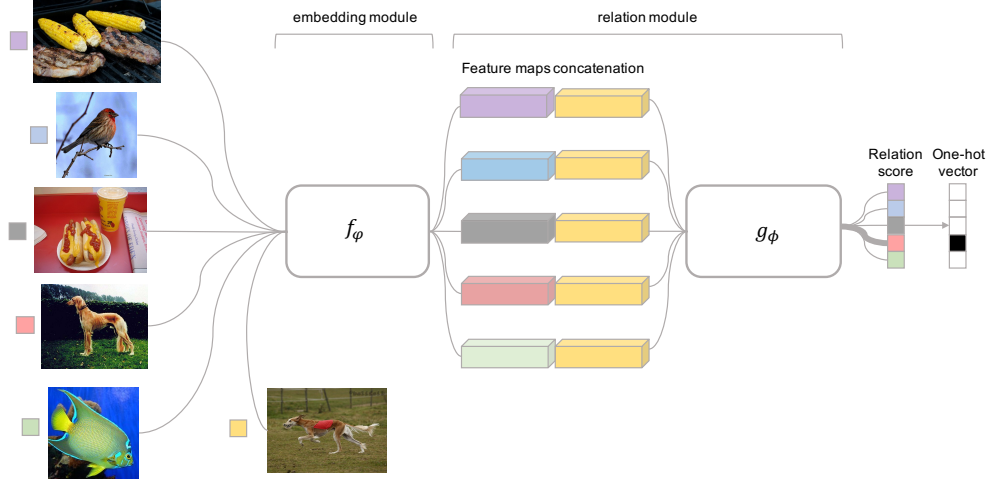


Figure 1: Relation Network architecture with a 5-way 1-shot 1-query example.

second heterogeneous embedding module f_{φ_2} besides the embedding module f_{φ_1} used for the image *query* set. Then the relation net g_ϕ is applied as before. Therefore, the relation score for each query input x_j will be:

$$r_{i,j} = g_\phi(\mathcal{C}(f_{\varphi_1}(v_c), f_{\varphi_2}(x_j))), \quad i = 1, 2, \dots, C \quad (3)$$

The objective function for zero-shot learning is the same as that for few-shot learning.

3.4. Network architecture

To demonstrate the generality of our approach, we instantiate our RN with multiple architectures.

Naive RN As most few-shot learning model utilise four convolutional blocks for embedding module [38, 35], we first introduce a *naive* version: Naive Relation Network (Naive RN), see Figure 2(b). For the embedding module, we use four convolutional blocks as per [38, 35]. Each convolutional blocks contains a 64-filter 3×3 convolution, batch normalisation and ReLU nonlinearity. The first two blocks contain a 2×2 max-pooling layer while the latter two do not. We do so because we need the output feature maps for further convolutional layers in the relation module. The relation module consists of two convolutional blocks and two fully-connected layers. Each of convolutional block is a 3×3 convolution with 64 filters followed by batch normalisation, ReLU non-linearity and 2×2 max-pooling. The output size of last max pooling $\mathcal{H} = 64$ and $\mathcal{H} = 64 * 3 * 3 = 576$ for Omniglot and *miniImageNet* respectively. The two fully-connected layers are 8 and 1 dimensional, respectively.

Deeper RN We also explore an architecture with *deeper* embedding and relation modules, Deeper Relation Network

(Deeper RN) (Figure 2(d)). The embedding module is composed of convolution blocks with 64-filter 7×7 convolution, batch normalisation, ReLU nonlinearity, 3×3 max-pooling and 3, 4 and 6 bottleneck blocks [14] with 64, 128, 256 filters respectively. The relation module contains 4 and 3 bottleneck blocks with 128, 64 filters – followed by average pooling and two fully-connected layers.

All fully-connected layers are ReLU except the output layer is Sigmoid in order to generate relation scores in a reasonable range for all versions of our network architecture. The zero-shot learning architecture is shown in Figure 2(e).

4. Experiments

We evaluate our approach on two related tasks: few-shot classification on Omniglot and *miniImageNet*, and zero-shot classification on Animals with Attributes (AwA) and Caltech-UCSD Birds-200-2011 (CUB). All the experiments are implemented based on PyTorch [1].

4.1. Few-shot Recognition

Settings Few-shot learning in all experiments uses Adam [17] with initial learning rate 10^{-3} , annealed by half for every 200,000 episodes. All our models are end-to-end trained from scratch with no additional dataset.

Baselines We compare against various state of the art baselines for few-shot recognition, including neural statistician [8], Matching Nets with and without fine-tuning [38], MANN [31], Siamese Nets with Memory [16], Convolutional Siamese Nets [18], MAML [10], Meta Nets [26], Prototypical Nets [35], Meta-Learner LSTM [28] and TCML [25].

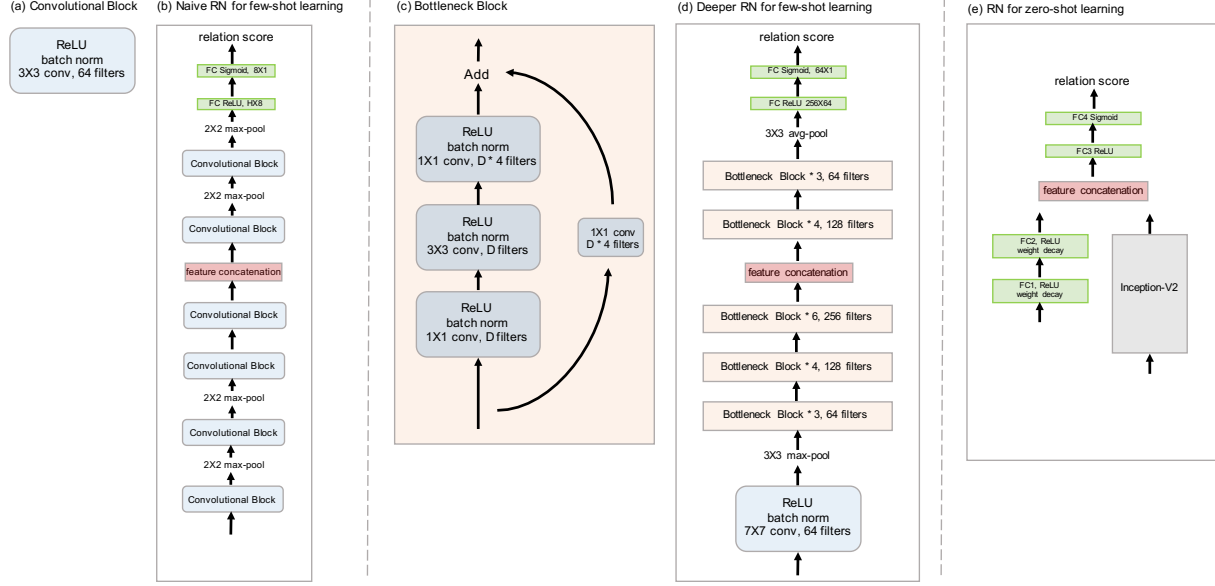


Figure 2: Relation Network architecture for few-shot learning: (b) naive version, (d) deeper version. Relation Network architecture for (e) zero-shot learning. These are composed of elements including (a) convolutional block, and (b) bottleneck block.

4.1.1 Omniglot

Dataset Omniglot [21] contains 1623 characters (classes) from 50 different alphabets. Each class contains 20 samples drawn by different people. Following [31, 38, 35], we augment new classes through 90° , 180° and 270° rotations of existing data and use 1200 original classes plus rotations for training and remaining 423 classes plus rotations for testing. All input images are resized to 28×28 .

Training Beside the K sample images, the **5-way 1-shot** contains 19 query images, the **5-way 5-shot** has 15 query images, the **20-way 1-shot** has 10 query images and the **20-way 5-shot** has 5 query images for each of the C sampled classes in each training episode. This means for example that there are $19 \times 5 + 1 \times 5 = 100$ images in one training episode/mini-batch for the 5-way 1-shot experiments.

Results Following [35], we computed few-shot classification accuracies on Omniglot by averaging over 1000 randomly generated episodes from the testing set. For the 1-shot and 5-shot experiments, we batch one and five query images per class respectively for evaluation during testing. The results are shown in Table 1. We achieved state-of-the-art performance under all experiments setting with higher averaged accuracies and lower standard deviations, except 5-way 5-shot where our model is 0.1% lower in accuracy than [10]. This is despite that many alternatives have significantly more complicated machinery [26, 8], or fine-tune on the target problem [10, 38] while we do not.

4.1.2 miniImageNet

Dataset The *miniImageNet* dataset, originally proposed by [38], consists of 60,000 colour images with 100 classes, each having 600 examples. We followed the split introduced by [28], with 64, 16, and 20 classes for training, validation and testing, respectively.

Training Following the standard setting adopted by most existing few-shot learning work, we conducted 5 way 1-shot and 5-shot classification. Beside the K sample images, the **5-way 1-shot** contains 15 query images, and the **5-way 5-shot** has 10 query images for each of the C sampled classes in each training episode. This means for example that there are $15 \times 5 + 1 \times 5 = 80$ images in one training episode/mini-batch for 5-way 1-shot experiments.

Architectures Almost all the alternatives use 4 convolutional blocks for the embedding network architecture. One exception is TCML [25] which uses a deeper embedding network. TCML [25] reported that with deeper embedding, MAML [10] obtained worse performance (30% on 1-shot *miniImageNet*), indicating that deeper networks do not *automatically* translate to better performance for few-shot classification. To compare with TCML [25] and also evaluate our approach with a deeper architecture, we evaluate both Deeper RN (Figure 2(d)) and Naive RN (Figure 2(b)). We resize input images to 224×224 for *deeper* version, 84×84 for *naive* version. Both networks are trained end-to-end from from scratch, with random initialisation, and no additional training set.

Results Following [35], we batch 15 query images per

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.

class in each episode for evaluation in both 1-shot and 5-shot scenarios and the few-shot classification accuracies are computed by averaging over 600 randomly generated episodes from the test set.

From Table 2, we can see that our *naive* version achieved state-of-the-art performance on all settings except 5-way 5-shot where intervals overlapped with [35]. However, the 1-shot result reported by prototypical networks [35] was trained on 30-way 15 queries per training episode, and 5-shot result was trained on 20-way 15 queries per training episode. When trained with 5-way 15 query per training episode, [35] only got $46.14 \pm 0.77\%$ for 1-shot evaluation, and $65.77 \pm 0.70\%$ for 5-shot evaluation, clearly weaker than ours. In contrast, all our models are trained on 5-way, 1 query for 1-shot and 5 queries for 5-shot per training episode, with much less training queries than [35].

Results are also shown when using *deeper* architecture (Figure 2(d)). We can see that: (i) Our approach is able to exploit the increased depth to obtain better performance compared to the basic 4-block version (as discussed, this is not automatic), (ii) Using a comparable architecture to TCML [25] we achieve state of the art performance on *miniImageNet*.

4.2. Zero-shot Recognition

Datasets Two widely used ZSL benchmarks are selected: **AwA** (Animals with Attributes) [22] consists of 30,745 images of 50 classes of animals. It has a fixed split for evaluation with 40 training classes and 10 test classes. **CUB** (Caltech-UCSD Birds-200-2011) [39] contains 11,788 images of 200 bird species. We use the same split as in [2] with 150 seen classes and 50 disjoint unseen classes.

Semantic representation For **AwA**, we use the continuous 85-dimension class-level attribute vector from [22], which has been used by all recent works. For **CUB**, a continuous 312-dimension class-level attribute vector is used.

Model	FT	5-way Acc.	
		1-shot	5-shot
MATCHING NETS [38]	N	43.56 \pm 0.84%	55.31 \pm 0.73%
META NETS [26]	N	49.21 \pm 0.96%	-
META-LEARN LSTM [28]	N	43.44 \pm 0.77%	60.60 \pm 0.71%
MAML [10]	Y	48.70 \pm 1.84%	63.11 \pm 0.92%
PROTOTYPICAL NETS [35]	N	49.42 \pm 0.78%	68.20 \pm 0.66%
RELATION NET (NAIVE)	N	51.38 \pm 0.82%	67.07 \pm 0.69%
TCML [25]	N	55.71 \pm 0.99%	68.88 \pm 0.92%
RELATION NET (DEEPER)	N	57.02 \pm 0.92%	71.07 \pm 0.69%

Table 2: Few-shot classification accuracies on *miniImageNet*. All accuracy results are averaged over 600 test episodes and are reported with 95% confidence intervals, same as [35]. For each task, the best-performing method is highlighted, along with any others whose confidence intervals overlap. ‘-’: not reported.

Implementation details Two different embedding modules are used for the two input modalities in zero-shot learning. Unless otherwise specified, we use Inception-V2 [37, 15] as the query image network in all ZSL experiments, taking the top pooling units as image embedding with dimension $D = 1,024$. This network is pre-trained on ILSVRC 2012 1K classification without fine-tuning, as in recent deep ZSL works [23, 29, 41]. A MLP network is used for embedding semantic attribute vectors. The size of hidden layer FC1 (Figure 2 (e)) is set to 800 and 900 for AwA and CUB respectively, and the output size FC2 is set to 1024 for both datasets. For the relation module, the 1024-D image and semantic embeddings are concatenated before being fed into MLPs with hidden layer FC3 size 800 and 900 for AwA and CUB, respectively. All ZSL models are trained with weight decay 10^{-5} in the embedding network. Learning rate is initialised to 10^{-5} with Adam [17] and then annealed by half every 200,000 iterations.

Conventional zero-shot learning The conventional eval-

uation setting for ZSL followed by the majority of prior work is to assume that the test data all comes from unseen classes [41]. We evaluate this setting first. We compare 15 alternative approaches in Table 3. With only the attribute vector used as the sample class embedding, our model achieves competitive result on AwA and state-of-the-art performance on the more challenging CUB dataset, outperforming the most related alternative prototypical networks [35] by a big margin. Note that only inductive methods are considered. Some recent methods [44, 12, 13] are transductive in that they use all test data at once for model training, which gives them a big advantage at the cost of making a very strong assumption that may not be met in practical applications, so we do not compare these here.

Generalised zero-shot learning Recently it has been argued that the traditional ZSL setting of solely testing-class data at runtime is unrealistic and too easy. The ‘generalised zero-shot learning’ setting which considers testing data as containing samples from both seen and unseen classes has therefore begun to gain increasing interest. We follow the setting of [6]. Specifically, we hold out 20% of the samples from the seen classes and mix them with the data samples from the unseen classes. The evaluation metric is now Area Under Seen-Unseen accuracy Curve (AUSUC), which measures how well a zero-shot learning method can trade-off between recognising data from seen classes and that of unseen classes [6]. Thus this metric is bounded by 1 and 0 above and below respectively. The results on AwA and CUB are presented in Table 4, comparing our model with six other alternatives. We can see that on AwA, our model significantly outperforms the competitors. While in CUB dataset, our method is only outperformed by SAE [19].

5. Further Analysis

5.1. Architectures

Loss function We used MSE loss (regression to one-hot vector), which is somewhat non-standard compared to the common cross-entropy. Table 5 shows a slightly inferior result is obtained when using cross entropy (CE) in place of MSE, in 1-shot and 5-shot *miniImageNet* experiments. In order to be the same as conventional classification, for the Cross-Entropy Version, we add another simple MLP after the relation module, then concatenate all relation scores fed into the MLP to obtain the softmax result. Then we use cross-entropy loss to train the whole network end-to-end.

How deep? The embedding module output feature maps which are concatenated and fed into our relation module. The efficacy of the embedding module and the richness of the feature map that they produce underpin the whole deep model. In this experiment, we vary the number of convolutional blocks in embedding module. Table 5 shows that

Model	F	SS	AwA	CUB
			10-way 0-shot	50-way 0-shot
SJE [2]	F_G	A	66.7	50.1
ESZSL [30]	F_G	A	76.3	47.2
SSE-RELU [42]	F_V	A	76.3	30.4
JLSE [43]	F_V	A	80.5	42.1
SYNC-STRUCT [5]	F_G	A	72.9	54.5
SEC-ML [4]	F_V	A	77.3	43.3
SAE [19]	N_G	A	84.7	61.4
PROTO. NETS [35]	N_G	A	-	54.6
DEVISE [11]	N_G	A/W	56.7/50.4	33.5
SOCHER <i>et al.</i> [36]	N_G	A/W	60.8/50.3	39.6
MTMDL [40]	N_G	A/W	63.7/55.3	32.3
BA <i>et al.</i> [23]	N_G	A/W	69.3/58.7	34.0
DS-SJE [29]	N_G	A/D	-	50.4/ 56.8
DEM [41]	N_G	A/W	86.7/78.8	58.3
RELATION NET	N_G	A	83.7	62.0

Table 3: Conventional zero-shot classification accuracy (%) comparison on AwA and CUB. SS: semantic space; A: attribute space; W: semantic word vector space; D: sentence description (only available for CUB). F: how the visual feature space is computed; For non-deep models: F_O if overfeat [33] is used; F_G for GoogLeNet [37]; and F_V for VGG net [34]. For neural network based methods, all use Inception-V2 (GoogLeNet with batch normalisation) [37, 15] as the CNN subnet, indicated as N_G .

Model	AwA	CUB
DAP [22]	0.366	0.194
IAP [22]	0.394	0.199
CONSE [27]	0.428	0.212
ESZSL [30]	0.449	0.243
SYNC-STRUCT [5]	0.583	0.356
SAE [19]	0.579	0.448
RELATION NET	0.650	0.371

Table 4: Performances measured in AUSUC of several methods for generalized Zero-shot learning on AwA and CUB. The higher the better (the upper bound is 1).

Model	1-shot	5-shot
2 CONV EMB. + CE	48.30 \pm 0.81%	62.10 \pm 0.72%
2 CONV EMB. + MSE	49.02 \pm 0.83%	63.80 \pm 0.70%
4 CONV EMB. + CE	49.60 \pm 0.80%	64.40 \pm 0.70%
4 CONV EMB. + MSE	51.38 \pm 0.82%	67.07 \pm 0.69%
DEEPER RN + MSE	57.02 \pm 0.92%	71.07 \pm 0.69%

Table 5: Effect of varying network depths and loss functions. Few-shot 5-way classification accuracies on *miniImageNet* averaged over 600 test episodes.

the performance on *miniImageNet* increases when the embedding module get deeper. Also the results of Deeper RN (Figure 2 (d)) in Table 5 show that our approach is able to exploit the increased depth to obtain better performance. All these models are trained from scratch.

5.2. Why does Relation Network Work?

Relationship to existing models Related prior few-shot work uses fixed pre-specified distance metrics such as Euclidean or cosine distance to perform classification [38, 35]. These studies can be seen as distance metric learning, but where all the learning occurs in the feature embedding, and a fixed metric is used given the learned embedding. Also related are conventional metric learning approaches [24, 7] that focus on learning a shallow (linear) Mahalanobis metric for a fixed feature representation. In contrast to prior work’s fixed metric or fixed features and shallow learned metric, Relation Network can be seen as both learning a deep embedding *and* learning a deep non-linear metric (similarity function)¹. These are mutually tuned end-to-end to support each other in few short learning.

Why might this be particularly useful? By using a flexible function approximator to learn similarity, we learn a good metric in a data driven way and do not have to manually choose the right metric (Euclidean, cosine, Mahalanobis). Fixed metrics like [38, 35] assume that features are solely compared element-wise, and the most related [35] assumes linear separability after the embedding. These are thus critically dependent on the efficacy of the learned embedding network, and hence limited by the extent to which the embedding networks generate inadequately discriminative representations. In contrast, by deep learning a non-linear similarity metric jointly with the embedding, Relation Network can better identify matching/mismatching pairs.

Visualisation To illustrate the previous point about adequacy of learned input embeddings, we show a synthetic example where existing approaches definitely fail and our Relation Network can succeed due to using a deep relation module. Assuming 2D query and sample input embeddings to a relation module, Fig. 3(a) shows the space of 2D sample inputs for a fixed 2D query input. Each sample input (pixel) is colored according to whether it matches the fixed query or not. This represents a case where the output of the embedding modules is not discriminative enough for trivial (Euclidean NN) comparison between query and sample set. In Fig. 3(c) we attempt to learn matching via a Mahalanobis metric learning relation module, and we can see the result is inadequate. In Fig. 3(d) we learn a further 2-hidden layer MLP embedding of query and sample inputs as well as the subsequent Mahalanobis metric, which is also not adequate. Only by learning the full deep relation module for similarity can we solve this problem in Fig. 3(b).

In a real problem the difficulty of comparing embeddings may not be this extreme, but it can still be challenging. We qualitatively illustrate the challenge of matching two exam-

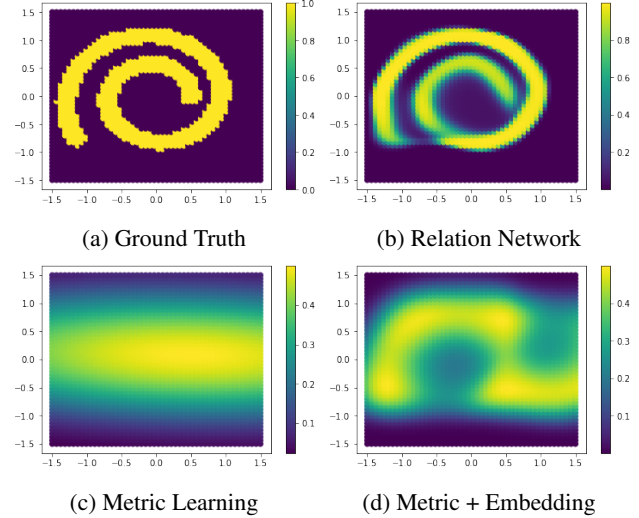


Figure 3: An example relation learnable by Relation Network and not by non-linear embedding + metric learning.

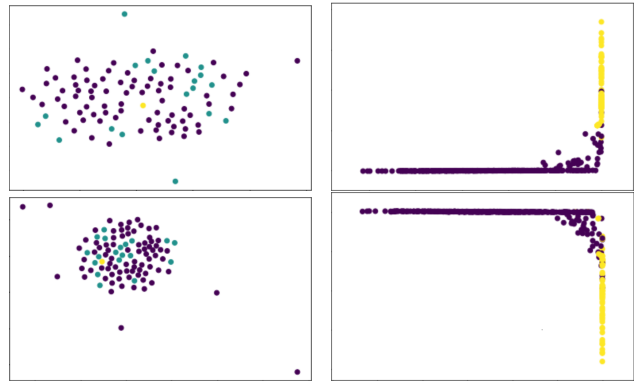


Figure 4: Example Omniglot few-shot problem visualisations. Left: Matched (cyan) and mismatched (magenta) sample embeddings for a given query (yellow) are not straightforward to differentiate. Right: Matched (yellow) and mismatched (magenta) relation module pair representations are linearly separable.

ple Omniglot query images (embeddings projected to 2D, Figure 4(left)) by showing an analogous plot of real sample images colored by match (cyan) or mismatch (magenta) to two example queries (yellow). Under standard assumptions [38, 35, 24, 7] the cyan matching samples should be nearest neighbours to the yellow query image with some metric (Euclidean, Cosine, Mahalanobis). But we can see that the match relation is more complex than this. In Figure 4(right), we instead plot the same two example queries in terms of a 2D PCA representation of each query-sample pair, as represented by the relation module’s penultimate layer. We can see that the relation network has mapped the data into a space where the (mis)matched pairs are linearly separable.

¹Our architecture does not guarantee the self-similarity and symmetry properties of a formal similarity function. But empirically we find these properties hold numerically for a trained Relation Network.

6. Conclusion

We proposed a simple method called the Relation Network for few-shot and zero-shot learning. Relation network learns an embedding and a deep non-linear distance metric for comparing query and sample items. Training the network end-to-end with episodic training tunes the embedding and distance metric for effective few-shot learning. This approach is far simpler and more efficient than recent few-shot meta-learning approaches, and produces state-of-the-art results. It further proves effective at both conventional and generalised zero-shot settings.

References

- [1] Pytorch. <https://github.com/pytorch/pytorch>. 4
- [2] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele. Evaluation of output embeddings for fine-grained image classification. In *CVPR*, 2015. 1, 3, 6, 7
- [3] L. Bertinetto, J. F. Henriques, J. Valmadre, P. H. S. Torr, and A. Vedaldi. Learning feed-forward one-shot learners. In *NIPS*, 2016. 1, 2
- [4] M. Bucher, S. Herbin, and F. Jurie. Improving semantic embedding consistency by metric learning for zero-shot classification. In *ECCV*, 2016. 7
- [5] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha. Synthesized classifiers for zero-shot learning. In *CVPR*, 2016. 7
- [6] W.-L. Chao, S. Changpinyo, B. Gong, and F. Sha. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *ECCV*, 2016. 7
- [7] D. Chen, X. Cao, L. Wang, F. Wen, and J. Sun. Bayesian face revisited: A joint formulation. In *ECCV*. Springer Berlin Heidelberg, 2012. 8
- [8] H. Edwards and A. Storkey. Towards a neural statistician. *ICLR*, 2017. 1, 4, 5, 6
- [9] L. Fei-Fei, R. Fergus, and P. Perona. One-shot learning of object categories. *TPAMI*, 2006. 1, 2
- [10] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 1, 2, 4, 5, 6
- [11] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A deep visual-semantic embedding model. In *NIPS*, 2013. 1, 3, 7
- [12] Y. Fu, T. M. Hospedales, T. Xiang, Z. Fu, and S. Gong. Transductive multi-view embedding for zero-shot recognition and annotation. In *ECCV*, 2014. 7
- [13] Y. Fu and L. Sigal. Semi-supervised vocabulary-informed learning. In *CVPR*, 2016. 7
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 4
- [15] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 6, 7
- [16] Ł. Kaiser, O. Nachum, A. Roy, and S. Bengio. Learning to remember rare events. *ICLR*, 2017. 1, 4, 6
- [17] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 4, 6
- [18] G. Koch, R. Zemel, and R. Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML Workshop*, 2015. 1, 2, 4, 6
- [19] E. Kodirov, T. Xiang, and S. Gong. Semantic autoencoder for zero-shot learning. In *CVPR*, 2017. 7
- [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 2
- [21] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum. One shot learning of simple visual concepts. In *CogSci*, 2011. 1, 2, 5
- [22] C. H. Lampert, H. Nickisch, and S. Harmeling. Attribute-based classification for zero-shot visual object categorization. *PAMI*, 2014. 1, 6, 7
- [23] J. Lei Ba, K. Swersky, S. Fidler, and R. Salakhutdinov. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *ICCV*, 2015. 1, 6, 7
- [24] T. Mensink, J. Verbeek, F. Perronnin, and G. Csurka. Metric learning for large scale image classification: Generalizing to new classes at near-zero cost. In *ECCV*, 2012. 8
- [25] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. Meta-learning with temporal convolutions. *arXiv preprint arXiv:1707.03141*, 2017. 4, 5, 6
- [26] T. Munkhdalai and H. Yu. Meta networks. In *ICML*, 2017. 1, 2, 4, 5, 6
- [27] M. Norouzi, T. Mikolov, S. Bengio, Y. Singer, J. Shlens, A. Frome, G. S. Corrado, and J. Dean. Zero-shot learning by convex combination of semantic embeddings. In *ICLR*, 2014. 7
- [28] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 1, 2, 4, 5, 6
- [29] S. Reed, Z. Akata, B. Schiele, and H. Lee. Learning deep representations of fine-grained visual descriptions. In *CVPR*, 2016. 6, 7
- [30] B. Romera-Paredes and P. Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015. 1, 7
- [31] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. Meta-learning with memory-augmented neural networks. In *ICML*, 2016. 1, 2, 4, 5, 6
- [32] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. A simple neural network module for relational reasoning. *arXiv preprint arXiv:1706.01427*, 2017. 2
- [33] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013. 7
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015. 1, 2, 7
- [35] J. Snell, K. Swersky, and R. S. Zemel. Prototypical networks for few-shot learning. In *NIPS*, 2017. 1, 2, 4, 5, 6, 7, 8
- [36] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, 2013. 7

- [37] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 6, 7
- [38] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *NIPS*, 2016. 1, 2, 3, 4, 5, 6, 8
- [39] C. Wah, S. Branson, P. Perona, and S. Belongie. Multiclass recognition and part localization with humans in the loop. In *ICCV*, 2011. 6
- [40] Y. Yang and T. M. Hospedales. A unified perspective on multi-domain and multi-task learning. In *ICLR*, 2015. 3, 7
- [41] L. Zhang, T. Xiang, and S. Gong. Learning a deep embedding model for zero-shot learning. In *CVPR*, 2017. 1, 6, 7
- [42] Z. Zhang and V. Saligrama. Zero-shot learning via semantic similarity embedding. In *ICCV*, 2015. 3, 7
- [43] Z. Zhang and V. Saligrama. Zero-shot learning via joint latent similarity embedding. In *CVPR*, 2016. 7
- [44] Z. Zhang and V. Saligrama. Zero-shot recognition via structured prediction. In *ECCV*, 2016. 7