# Siamese Neural Networks for One-shot Image Recognition

*Akriti Pandey[1], Rani Gupta [2], Rupali Dubey[3] Mrs. Shikha Agrawal [4]*

[1-4]Department of Information Technology, RKGIT, UP, India

## *Abstract*

Learning excellent features for machine learning applications can be computationally expensive, and it can be difficult in situations when there is limited data. The one-shot learning setting is a classic example of this, in which we must accurately make predictions based on only one sample of each new class.

We look at a strategy for learning Siamese neural networks that uses a unique structure to automatically priorities similarity between inputs in this research. After tuning a network, we can use powerful discriminative features to extend the network's predictive capacity not merely to fresh data, but to totally new classes from unknown distributions. We can generate strong results using a convolutional design that outperforms the competition.

*Keywords:* Siamese neural networks, machine learning, One-shot learning, image recognition, Omniglot data set

---

## INTRODUCTION

## 1. INTRODUCTION

Humans have a remarkable capacity for learning and recognizing new patterns. We have shown that when people are exposed with stimuli, they appear to be able to swiftly grasp new concepts and subsequently recognize variations on these concepts in future perceptions. This technique is applicable to more than only generalizing to unseen images of a concept class. When shown an image of a specific species, such as a tiger, a person may not only identify new photographs of tigers using that image, but also make accurate estimates regarding images of new species. For example, we might be able to say that a lion resembles a tiger more than another mammal. As a result, this form of inference can be extremely useful.

In a range of applications, such as web search, spam detection, caption generation, and speech and picture recognition, machine learning has been effectively employed to reach state-of-the-art performance.
However, when these algorithms are pushed to generate predictions about data for which there is little supervised knowledge, they frequently fail. We want to be able to generalize to these new categories without requiring lengthy retraining, which could be costly or difficult owing to a lack of data, or in an online prediction context like web retrieval.

Classification under the constraint of just seeing a single example of each potential class before generating a prediction about a test instance is a particularly intriguing task. This is known as one-shot learning, and it is the main emphasis of the model we provide in this paper. This is distinct from zero-shot learning, in which the model is unable to see any examples from the target classes.

One-shot learning can be directly addressed by creating domain-specific features or inference processes with highly discriminative properties for the target task. As a result, systems that include these strategies succeed in similar situations but fall short of providing strong solutions that can be applied to a variety of issues. We describe a novel strategy in this research that reduces input structural assumptions while automatically collecting characteristics that allow the model to generalize well from a small number of samples. We build on the deep learning framework, which employs many layers of nonlinearities to capture transformation invariances in the input space, often by constructing a model with many parameters and then analysing a significant amount of data.

## 2. LITERATURE REVIEW AND OBJECTIVE

Overall, research on one-shot learning algorithms is still in its early stages, and the machine learning community has paid little attention to it. There are, however, a few crucial lines of work that come before this study. Although a few scholars worked on one-shot learning in the 1980s and 1990s, the foundational work on applying machine learning to one-shot learning was published in the early 2000s. Li Fei-Fei et al. created a variational Bayesian framework for one-shot picture classification based on the idea that previously learnt classes can be used to predict future ones when only a few samples are available. This data can be successfully incorporated into the prior, which is updated as new classes are discovered, and then merged with the likelihood to produce a new posterior distribution for the class The authors created a generative Constellation model with appearance and shape components that detects discrete feature space regions..

Lake et al., more recently, took a cognitive science approach to the problem of one-shot learning, tackling one-shot learning for character recognition with a method called Hierarchical Bayesian Program Learning (HBPL).

One-shot learning approaches in computer vision are usually divided into two categories: feature learning and metric learning. Wan et al. developed an intricate SIFT-based hierarchical feature extraction approach for which the generated features are passed to a nearest-neighbor classifier in their work on gesture detection difficulties.

Instead of explicitly addressing one-shot learning, the authors focus on how to "borrow" examples from other classes in the training set. This approach might be effective for data sets with few examples for some classes, since it allows for a flexible and ongoing integration of inter-class information into the model. The authors present a new loss function that regularises an adjustable weight vector representing a soft measure of how much each category should borrow from the current training exemplar.

## 3. MATERIALS AND METHODS

Although the fundamental approach can be duplicated for nearly any modality, we focus on character recognition.

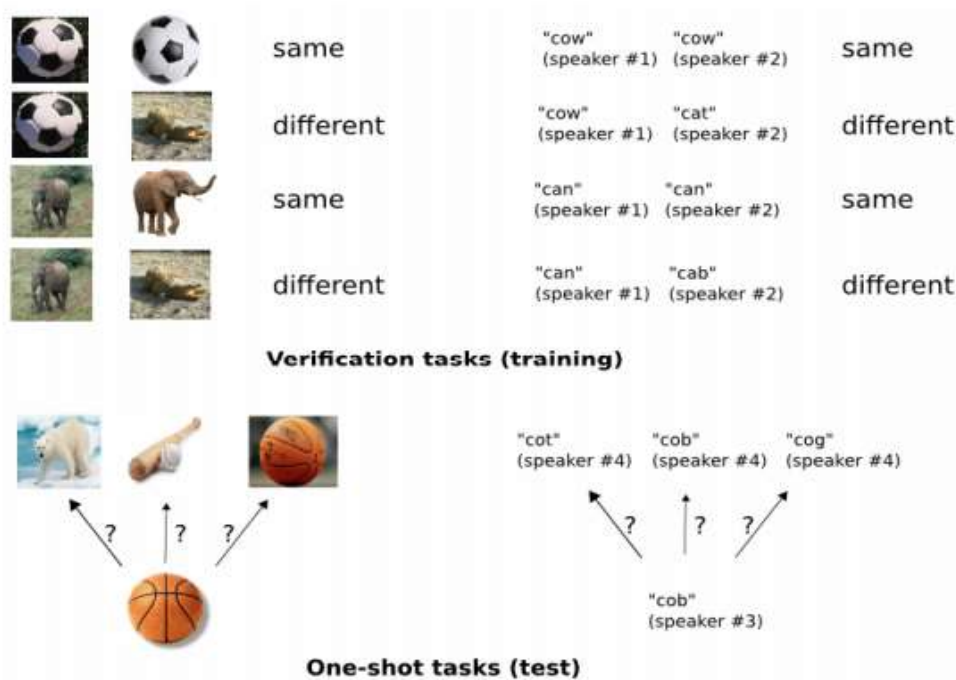This is our overall plan :-

1) Train a model to distinguish between a set of similar and dissimilar pairs.

2) For verification, generalise to evaluate new categories based on learnt feature mappings.

For this domain, we use large siamese convolutional neural networks, which:

a) can learn generic image features useful for making predictions about unknown class distributions even when there are few examples from these new distributions;

b) can be easily trained using standard optimization techniques on pairs sampled from the source data; and

c) provide a competitive approach that does not rely on domain-specific knowledge by instead exploiting.

We want to construct a neural network that can discriminate between the class-identity of image pairings, which is the typical verification challenge for image recognition, before developing a model for one-shot image classification. This network will predict whether two photos from the same alphabet (or potentially distinct alphabets, however we do not explore this situation) portray the same character.

**The Omniglot data set-** Brenden Lake and his colleagues at MIT used Amazon's Mechanical Turk to collect the Omniglot data set in order to create a common benchmark for learning from few instances in the handwritten character recognition domain. Omniglot includes samples from 50 alphabets, including well-known worldwide languages like Latin and Korean as well as lesser-known regional dialects. There are also some fictional character sets like Aurek-Besh and Klingon.

Each alphabet has a different number of letters, ranging from roughly 15 to more than 40. Each of the 20 drawers produces all of the characters in these alphabets once. The overall data set contains only a few samples for each potential letter class; as a result, the authors refer to it as a "MNIST transposition," in which the number of classes significantly outnumbers the number of training cases.

Lake divided the data into two sets: a 40-letter backdrop set and a 10-letter evaluation set. We keep these two words to distinguish ourselves from the standard training, validation, and test sets that can be constructed from the background set to tune models for verification. The background set is used to learn hyperparameters and feature mappings in order to build a model. The assessment set, on the other hand, is only utilised to assess one-shot classification performance. The phrases "background set" and "evaluation set" are used throughout this study in the same way they are in Lake's work.

One-Shot Learning Task-

Lake created a 20-way within-alphabet classification task to empirically evaluate one-shot learning performance, in which an alphabet is initially chosen from among those reserved for the evaluation set.

Omniglot's characters are 105x105 binary-valued images produced by hand on an online canvas. The most recent version of the data set includes all of the alphabets, as well as the exact one-shot trials utilised in Lake's original work and the stroke trajectories collected during the data set's development. The stroke trajectories were gathered alongside the composite images, allowing models trained on Omniglot to add temporal and structural information.

The Omniglot dataset contains a variety of different images from alphabets across the world. Example of a 20-way one-shot classification task using the Omniglot dataset-

The lone test image is displayed atop a grid of 20 photos that represent the various unseen classes from which we might select the test image. These 20 photos are the sole instances of each of those classes that we are aware of. in addition to twenty characters chosen at random.

Two of the twenty drawers are also chosen from the evaluation drawers pool. The twenty characters are then sampled by these two drawers. Each of the first drawer's characters is designated as a test picture, and it is compared to each of the second drawer's twenty characters, with the purpose of identifying the class matching to the test image from among all of the second drawer's characters.

For each of the ten evaluation alphabets, this process is done twice, for a total of 40 one-shot learning trials. There are 400 one-shot learning trials in total, from which the standard classification accuracy is computed.

We used the same collection of one-shot learning challenges as in [16] to replicate this technique. We also created software that can generate new one-shot tasks from any data set. This allows us to track one-shot learning performance while optimising for the verification task on our validation set.

## 4. RESULTS AND DISCUSSION

The Omniglot data set has only a few samples for each potential letter class; as a result, the authors describe to it as a "MNIST transposition," in which the number of classes significantly outnumbers the number of training cases (Lake et al., 2013). We thought it would be interesting to see how well an Omniglot-trained model might generalise to MNIST, where the 10 digits in MNIST are treated as an alphabet and a 10-way oneshot classification task is evaluated. On the MNIST test set, we used a similar technique to Omniglot, generating 400 one-shot trials but excluding any fine tuning on the training set. All 28x28 photos were down sampled to 35x35 and then fed to a reduced version of our model that had been trained on Omniglot 35x35 images that had been down sampled by a factor of three. On this job, we also looked at the nearest-neighbour baseline.

The nearest neighbour baseline performs similarly to Omniglot, whereas the convolutional network's performance suffers a more severe reduction. However, without any training on MNIST, we can still achieve adequate generalisation from the features learnt on Ominglot.

## 5. CONCLUSIONS

We've described a method for doing one-shot classification that involves learning deep convolutional siamese neural networks for verification first. We presented fresh findings comparing our networks' performance to that of an existing state-of-the-art classifier designed for the Omniglot data set. Our networks exceed all existing baselines by a large margin and come close to the best results obtained by prior authors. We've argued that these networks' outstanding performance on this test shows not only that our metric learning approach can achieve human-level accuracy, but also that it should be applied to one-shot learning tasks in other areas, particularly picture classification.

## REFERENCES

[1] Jane Bromley, James W. Bentz, Leon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Edouard Sackinger, and Roopak Shah. A siamese time delay neural network is used to verify signatures. 7

(04):669–688, 1993, International Journal of Pattern Recognition and Artificial Intelligence.C. D. Rakopoulos, and E. G. Giakoumis, Second-Law Analyses Applied to Internal Combustion Engine Operation, Progress in Energy and Combustion Science, 32(1), 2006, pp. 2-47.

[2]   Yoshua Bengio, Yoshua Bengio, Yoshua Bengio, Yoshua For AI, deep architectures are being learned. 2(1):1–127, 2009. Foundations and Trends in Machine Learning.

[3]   Di Wu, Fan Zhu, and Ling Shao. Learning gesture recognition from RGBD photos in one shot. The 2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 7–12. IEEE, 2012.

[4]   Li Fe-Fei, Robert Fergus, and Pietro Perona. Unsupervised one-shot learning of object categories using a Bayesian technique. In the year 2003, a paper was published in the journal Computer Vision. Proceedings. Pages. 1134–1141 in Ninth IEEE International Conference on. 2003, IEEE

[5]   Brenden M. Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua B. Tenenbaum. Simple visual concepts can be learned in a single shot. In Proceedings of the Cognitive Science Society's 33rd Annual Conference, volume 172, 2011..