

FuCiTNet: Improving the generalization of deep learning networks by the fusion of learned class-inherent transformations

Manuel Rey-Area^a, Emilio Guirado^b, Siham Tabik^c, Javier Ruiz-Hidalgo^d

^a*atlanTTic Research Center for Telecommunication Technologies, University of Vigo, Galicia, Spain*

^b*Multidisciplinary Institute for Environment Studies "Ramon Margalef" University of Alicante, Alicante, Spain*

^c*Andalusian Research Institute in Data Science and Computational Intelligence University of Granada 18071, Spain*

^d*Department of Signal Theory and Communications, Universitat Politècnica de Catalunya, Barcelona, Catalonia, Spain*

Abstract

It is widely known that very small datasets produce overfitting in Deep Neural Networks (DNNs), i.e., the network becomes highly biased to the data it has been trained on. This issue is often alleviated using transfer learning, regularization techniques and/or data augmentation. This work presents a new approach, independent but complementary to the previous mentioned techniques, for improving the generalization of DNNs on very small datasets in which the involved classes share many visual features. The proposed model, called FuCiTNet (Fusion Class inherent Transformations Network), inspired by GANs, creates as many generators as classes in the problem. Each generator, k , learns the transformations that bring the input image into the k -class domain. We introduce a classification loss in the generators to drive the learning of specific k -class transformations. Our experiments demonstrate that the proposed transformations improve the generalization of the classification model in three diverse datasets.

1. Introduction

It is well-known that building robust and efficient supervised deep learning networks requires large amounts of quality data, especially when the involved classes share many visual features. Recent technological advances in camera sensors and the potentially unlimited data provided by internet have helped gathering large volumes of images [1, 2, 3]. However, labelling such amounts of data is still manual and costly. In consequence, a large number of problems has still to deal with very small labeled datasets and hence, the resulting classification models are usually unable to generalise correctly to new unseen examples, this problem is known as overfitting.

In general, the issue of overfitting in supervised networks is addressed either from the model or data point-of-view. In the former, transfer learning or diverse regularization techniques are employed. In the later, increasing the volume of the training set using data augmentation strategies [4, 5, 6] are considered. For highly sensitive purposes, data augmentation strategies, if not chosen correctly, can potentially change meaningful information resulting in ill-posed training data. Instead of manually selecting data augmentation techniques, recent works showed that learning these transformations from data can lead to significant improvements in the generalization of the models.

Some approaches train the generator of a GAN (Generative Adversarial Network) to learn the suitable augmenting techniques from scratch [7, 8] while others train the generator to learn finding the optimal set of augmenting techniques from an initial space of data augmentation strategies [9]. The downside of the aforementioned approaches is that the selected transformations are applied to all samples in a training set without taking into account the particularities of each class.

The present work proposes FuCiTNet approach for learning class-inherent transformations in each class within the dataset under evaluation. FuCiTNet is inspired by GANs, it creates a number, N , of generators equal to the number of classes. Each generator, k , with $k \in \{1, \dots, N\}$, will be entrusted to learn the features of a specific k -class space. When a sample is fed into the system, it is broadcasted to every generator producing N transformed images, each of which is fed to the classifier which predicts a label with certain error. The error is transferred back to the entrusted generator (specified by the input's label ground-truth) indicating the amount of change the class transformation must be altered to meet the classifier requirements. The final prediction of the trained classification model will be calculated based on the fusion of the N different output predictions.

The contributions of this work can be summarized as follows:

- We propose class-inherent transformation generators for improving the generalization capacity of image classification models, especially appropriate for problems in which the involved classes share many visual features. Our approach, FuCiTNet, creates as many generators as the number of classes N . Each generator, $k \in \{1, \dots, N\}$, learns the inherent transformations that bring the input image from its space to the k -class space. We introduce a classification loss in the generators to drive the learning of specific k -class transformations.
- The final prediction of the classification model is calculated as a fusion of the N output scores. The source code of FuCiTNet will be available in Github after acceptance
- Our experiments demonstrate that class-inherent transformations produce a clearer discrimination for the classifier yielding better generalisation performance in three small datasets from two different fields.

This paper is organized as follows. A summary of the most related works to ours are reviewed in Section 2. A description of FuCiTNet model is provided in Section 3. Experimental framework is provided in Section 4. Results and analysis are given in Section 5 and finally conclusions and future work in Section 6.

2. Related Work

Improving the performance of supervised deep neural networks in image classification is still ongoing research. To reach high accuracies the model needs to generalise robustly to unseen cases to eventually avoid overfitting. This is addressed using several approaches.

The most popular approach is data augmentation. It was firstly introduced by Y. LeCun *et al.* [10] by randomly applying these transformations to the training dataset: shearing, random cropping and rotations. With the revolution of CNNs [4], novel transformations appeared, such as horizontal flipping, changes in intensity, zooming, noise injection etc. With data augmentation, most CNN-based classifiers reduce overfitting and increase their classification accuracy.

Dropout, a well known regularization technique, is also used for improving generalization. It was first introduced by Srivastava *et al.* [11]. The key idea is to randomly drop units (along with their connections) from the neural network throughout the training process. This prevents units from co-adapting too much to the data. It significantly reduces overfitting and gives major improvements over other regularization methods.

The emergence of GANs [12] has led to promising results in image classification. They were useful for generating synthetic samples to increase small datasets where the number of samples per class was low eventually introducing variability for the generalisation of the classification models. The downside of GANs is that its latent space converges if there exists a fair good amount of images to train. Generating synthetic samples in small and very small dataset is still an open issue.

For small and very small datasets, the problem of overfitting is even greater. Inspired by GANs, the authors in [7] designed an augmentation network that learns the set of augmentations that best improve the classifier performance. The classifier tells the generator network which configuration of image transformations prefers when distinguishing samples from different classes. In the same direction, the authors in [13, 9] address a similar problem using Reinforcement Learning and adversarial strategy respectively. The former chooses from a list of potential transformations which one is the best suited through augmentation policies. The latter combined the list of transformations to synthesize a total image transformation. The present work is different to all the previously cited works in that it proposes a new approach for learning transformations inherent to each particular class based on the classifier requirements eventually forcing the classes to be as distinguishable as possible from each other.

3. FuCiTNet approach

Inspired by GANs, FuCiTNet learns class-inherent transformations. Our aim is to build a generator that improves the discrimination capacity between the different classes. GANs use the so called adversarial loss which optimizes a min-max problem. The generator, G , tries to minimize the following function while the discriminator, D , tries to maximize it:

$$\min_G \max_D \mathbb{E}_{X \sim p_{data(X)}} [\log D(X)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

Where $D(X)$ is the discriminator’s estimate of the probability that real data instance X is real. $\mathbb{E}_{X \sim p_{data(X)}}$ is the expected value over all real data instances. $G(z)$ is the generator’s output given a noise z . $D(G(z))$ is the discriminator’s estimate of the probability that a fake instance is real. $\mathbb{E}_{z \sim p_z(z)}$ is the expected value over all random inputs to the generator. The formula derives from the cross-entropy between the real and generated distributions.

In this work, the adversarial concept is slightly changed. Instead of using a discriminator that focuses on maximizing the inter-class probability of the image belonging to the data or latent distribution, we use a classifier that minimizes the loss of $G(X)$ of belonging to a particular dataset class. As the images produced by the generator do not need to be similar to its ground truth, we restrict the generator latent space in a different way by replacing $D(X)$ with the classification loss of $G(X)$, which will hallucinate samples with specific enhanced features to meet the classifier requirements improving the classification accuracy.

Instead of building one single generator G , we create an array of N generators, G_k with $k \in \{1, \dots, N\}$. Where N is equal the number of object classes. Each generator $G_k(X) = X'_k$ is in charge of learning the inherent features of that specific k-class until X' becomes part of the k-class space. In other words, each generator G_k will learn the transformations that map the input image X from its own i-class domain to the k-class domain, with $i \in \{1, \dots, N\}$. The flowchart diagram of FuCiTNet is depicted in Figure 1.

The architecture of our generators G_k , with $k \in \{1, \dots, N\}$, consists of 5 identical residual blocks [14]. Each block has two convolutional layers with 3×3 kernels and 64 feature maps followed by batch-normalization layers [15] and ParametricReLU [16] as activation function. The last residual block is followed by a final convolutional layer which reduces the output image channels to 3 to match the input’s dimensions.

Our classifier is a ResNet-18 which consists of an initial convolutional layer with 7×7 kernels and 64 feature maps followed by a 3×3 max pool layer. Then, 4 blocks of two convolutional layers with 3×3 kernels with 64, 128, 256 and 512 feature maps respectively followed by a 7×7 average pooling and a fully connected layer which outputs a N element vector. ReLU is used as the activation function [17].

The classifier uses a standard cross entropy loss (l_{CE}). Given a batch size B of input images with their respective ground truth labels, $\{X^{(j)}, y^{(j)} \mid j =$

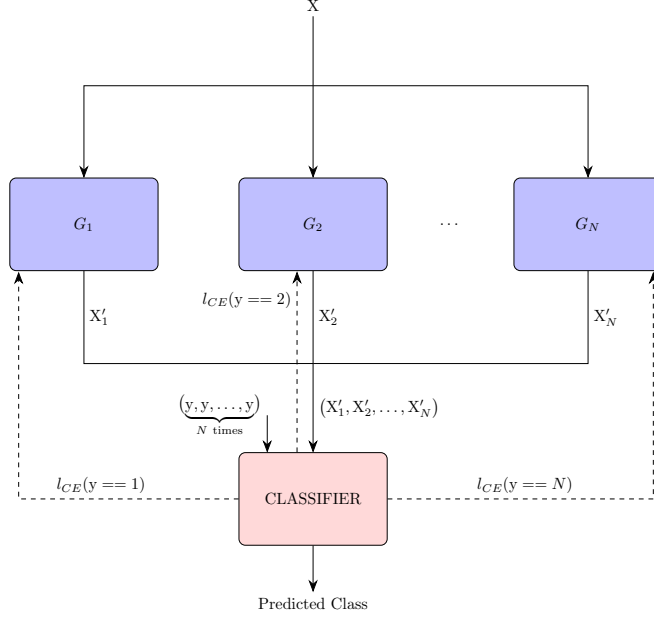


Figure 1: Flowchart of FuCiTNet during training. The input image X is broadcasted to every generator. Each generator produces a transformed image X' . The classifier computes the cross entropy loss which is transferred back to the generator commissioned to enhance features of the class given by the input's groundtruth label y .

$1, \dots, B\}$:

$$l_{CE} = -\frac{1}{B} \sum_{j=1}^B \sum_{k=1}^N \sum_{i=1}^N y^{(j)}[i] \cdot \log(\hat{y}_k^{(j)}[i]) \quad (2)$$

The true label distribution for the j^{th} image is depicted by $y^{(j)}$ and $\hat{y}_k^{(j)}$ indicates the predicted label distribution for the transformed image $X_k^{(j)}$.

All generators use the same loss inspired by [18]. For the k^{th} generator with $k \in \{1, \dots, N\}$ the respective loss \mathcal{L}_{gen_k} is indicated by Formula (3). It consists of a multiple term loss constituted by a pixel-wise MSE term, a perception MSE term and the classifier loss. Adding this classifier loss to the generators is one of the main novelty of our method. The classification loss is added to the generator loss with a weighted factor λ indicating how much the generator must change its outcome to suit the classifier. It is worth noting that, for each input image, the classification loss is only transferred to a particular generator indicated by the input's ground-truth label.

$$\mathcal{L}_{gen_k} = \underbrace{l_{MSE} + 0.006 \cdot l_{Perceptual}}_{\text{Similarity term: } \mathcal{L}_{sim}} + \lambda \cdot l_{CE}(y == k) \quad (3)$$

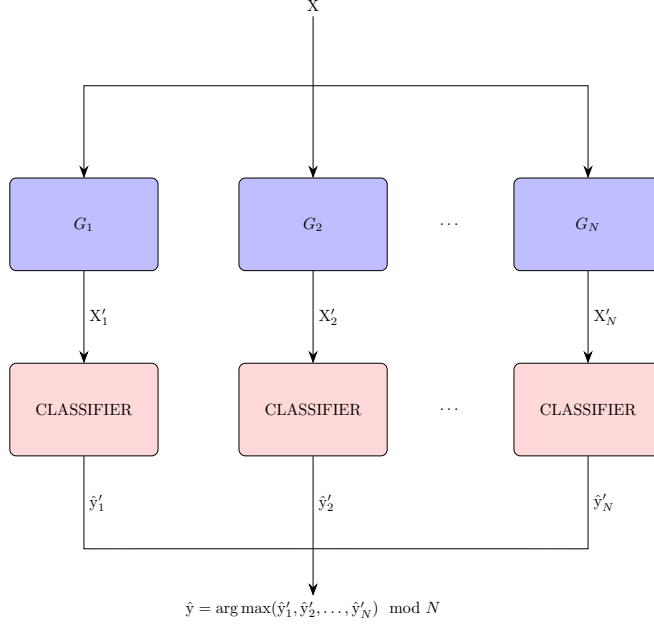


Figure 2: Flowchart of FuCiTNet in inference.

The similarity term (\mathcal{L}_{sim}) keeps the generator from changing the image too much while the classifier loss drives the output away from the input and close to the k-class feature space.

The pixel-wise MSE is obtained by performing a regular L2-norm between each pixel in the input image and the generated image.

$$l_{MSE} = \frac{1}{WH} \sum_{x=1}^W \sum_{y=1}^H (X_{x,y} - X'_{x,y})^2 \quad (4)$$

The perceptual MSE assesses similarity between X and X' by feeding each of them into a pretrained VGG-16 [19]. The euclidean distance between the j^{th} VGG-16 feature maps (ϕ_n) defines the perceptual loss.

$$l_{perceptual} = \frac{1}{W_j H_j} \sum_{x=1}^{W_j} \sum_{y=1}^{H_j} (\phi_j(X_{x,y}) - \phi_j(X'_{x,y}))^2 \quad (5)$$

The cross entropy loss coming from the classifier is added to a particular generator loss in the generator array specified by the input's ground-truth label y with a weighted factor λ controlling the impact of the classifier within the generator:

$$l_{CE}(y == k) = -\frac{1}{B} \sum_{j=1}^B \sum_{k=1}^N y^{(j)}[k] \cdot \log(\hat{y}'_k^{(j)}[k]) \quad (6)$$

where B is the batch size, N the number of classes and $y^{(j)}[k]$ is the element in the true label distribution of the j^{th} image belonging to class k . Likewise, $\hat{y}_k^{(j)}[k]$ depicts the classifier predictions for the transformed image $X_k^{(j)}$ with $k \in \{1, \dots, N\}$ of belonging to the class k .

In this manner, each generator in the array is able to learn class-inherent features from a specific class eventually, accentuating the differentiation among classes.

In inference time the general flowchart of the system is modified as shown in Figure 2. The final class prediction is given by concatenating each output distribution \hat{y}' for every X' in the logit domain, taking the arg max and computing the modulo with the amount of classes N as indicated in Formula (7)

$$\hat{y} = \arg \max(\hat{y}'_1, \hat{y}'_2, \dots, \hat{y}'_N) \mod N \quad (7)$$

4. Experimental framework

To assess the capacity of FuCiTNet in increasing the discrimination among classes, we created three datasets from two different fields using object-classes that are frequently confused by the best performing classification model. In particular, from Tiny ImageNet [20], we created two datasets made of two and three classes respectively, cat-vs-dog and cat-vs-dog-vs-goldfish. From NWPU-RESISC45 remote sensing dataset made of aerial ortho-images [21], we created church-vs-palace dataset which represent one of the most similar/confusing pair of classes in this dataset. To further increase the complexity of this dataset, we downsampled the church and palace images to 64×64 pixels. A brief description of the three datasets is provided in Table 1.

Dataset	# classes	# pixels/image	Train	Test	Total
cat-vs-dog	two	64×64	800	200	1000
cat-vs-dog-vs-goldfish	three	64×64	1200	300	1500
church-vs-palace	two	64×64	1120	280	1400

Table 1: Description of the three evaluated datasets together with the used data distribution 80%-20%.

The N generator networks were initialized randomly while the classifier, ResNet-18, was initialized using the pretrained weights on ImageNet [1]. As optimization algorithm, we used Adam [22] in both, generator and classifier with $\beta_1 = 0.9$. The adopted learning rates for generators and classifier are different, for the generators, we used a value of 10^{-4} and 10^{-3} for ResNet-18. The reason why generators have a lower learning rate is because they need to be subtle when generating the transformed image and avoid the classifier from falling behind on capturing the rate of change in appearance. For the ResNet-18 we used learning rate decay of 0.1 each 5 epochs. We have also used weight decay of 10^{-4} to avoid overfitting. We applied early stopping monitoring based on the validation loss with a patience of 10 epochs.

The weighted factor λ in the generator loss in Eq. (3) is a hyperparameter to be tuned. We evaluated a space of 13 values: [1, 0.5, 0.1, 0.075, 0.05, 0.025, 0.01, 0.0075, 0.005, 0.0025, 0.001, 0.00075, 0.0005]. We assessed the effect from a high contribution in the loss towards a softer impact. For each value of λ the system was trained throughout 100 epochs with a batch size of 32. We alternate updates between the generator and classifier, i.e., for each batch we update first the classifier then we update the generator. We evaluate FuCiTNet on the validation set after each epoch. The chosen weights for both networks are the ones which minimize the classification loss in this set.

For a fair comparison, we compare the results of FuCiTNet with the two most related approaches [7] and [9]. As [9] provides the source code, we analyze and show the results on the three considered datasets, cat-vs-dog, cat-vs-dog-vs-goldfish and church-vs-palace, following the same experimental protocol we used in the rest of experiments. However, as [7] does not provide the source code and considered only binary classes, we included only the results reported in the paper using the same experimental protocol on dog-vs-cat dataset. The results of this approach on cat-vs-dog-vs-goldfish and church-vs-palace are not available to us. Both approaches from [7] and [9] used ResNet CNN architecture. In addition, we also compare our results to the results of the best classification model obtained based on the same network architecture and by manually selecting the set of optimizations that reaches the highest performance.

In all the experiment, we used 3-fold cross validation following a 80:20 hold out data distribution as depicted in Table 1. All the experiments were executed on a NVIDIA Titan Xp. All the implementations were performed using PyTorch DL framework [23].

5. Results and analysis

This section presents, compares and analyzes the quantitative and qualitative results of FuCiTNet with the state-of-the-art methods, [7] and [9], and with the best classification model.

Setup	Accuracy	Cat mean confidence	Dog mean confidence
None	0.872	—	—
FT	0.865	—	—
Data aug	0.892	1.673	2.520
Data aug, FT	0.888	—	—
[7]	0.770	—	—
[9]	0.800	—	—
FuCiTNet, None, $\lambda = 0.01$	0.870	—	—
FuCiTNet, FT, $\lambda = 0.025$	0.880	—	—
FuCiTNet, Data aug, $\lambda = 0.0075$	0.883	—	—
FuCiTNet, Data aug, FT, $\lambda = 0.05$	0.912	2.734	2.101

Table 2: Performance of ResNet-18 classification model without (row: 1 to 4) and with (row: 7 to 10) FuCiTNet, using different configurations, on cat-vs-dog dataset. Data aug consists of random horizontal flipping, random rotation and random affine. The accuracy of the state-of-the-art approaches, [7] and [9], is shown in row 5 and 6 respectively.

Predicted/Actual	Cat	Dog
Cat	86	14
Dog	11	89

(a) Confusion matrix for the setup: Data aug

Predicted/Actual	Cat	Dog
Cat	97	3
Dog	10	90

(b) Confusion matrix for the setup: FuCiTNet, Data aug, FT, $\lambda = 0.05$

Table 3: Cat-vs-dog confusion matrices for fold #1 in 3FCV for the best reference model(a) and FuCiTNet (b)

Setup	Accuracy	Church mean confidence	Palace mean confidence
None	0.774	—	—
No data aug, FT	0.769	—	—
Data aug	0.783	0.769	1.452
Data aug, FT	0.779	—	—
[9]	0.779	—	—
FuCiTNet, None, $\lambda = 0.05$	0.751	—	—
FuCiTNet, FT, $\lambda = 0.0025$	0.767	—	—
FuCiTNet, Data aug, $\lambda = 0.01$	0.777	—	—
FuCiTNet, Data aug, FT, $\lambda = 0.0025$	0.795	1.868	1.580

Table 4: Performance of ResNet-18 based classification model without (row: 1 to 4) and with (row: 6 to 9) FuCiTNet, using different configurations, on the church-vs-palace dataset. Data augmentation consists of random horizontal flipping, random rotation and random affine. The accuracy of the state-of-the-art approach [9] is shown in row 5.

Predicted/Actual	Church	Palace
Church	114	26
Palace	30	110

(a) Confusion matrix for the setup: Data aug

Predicted/Actual	Church	Palace
Church	114	26
Palace	27	113

(b) Confusion matrix for the setup: FuCiTNet, Data aug, FT, $\lambda = 0.0025$

Table 5: Church-vs-palace confusion matrices for fold #1 in 3FCV for the best reference model(a) and FuCiTNet (b)

Setup	Accuracy	Cat mean confidence	Dog mean confidence	Goldfish mean confidence
None	0.898	—	—	—
FT	0.909	—	—	—
Data aug	0.915	7.796	4.778	9.617
Data aug, FT	0.911	—	—	—
[9]	0.790	—	—	—
FuCiTNet, None, $\lambda = 0.0075$	0.887	—	—	—
FuCiTNet, FT, $\lambda = 0.0075$	0.902	—	—	—
FuCiTNet, Data aug, $\lambda = 0.005$	0.900	—	—	—
FuCiTNet, Data aug, FT, $\lambda = 0.005$	0.920	3.484	3.248	6.613

Table 6: Performance of ResNet-18 based classification model without (row: 1 to 4) and with (row: 6 to 9) FuCiTNet, using different configurations, on cat-vs-dog-vs-goldfish dataset. Data augmentation consists of random horizontal flipping, random rotation and random affine. The accuracy of the state-of-the-art approach [9] is shown in row 5.

Predicted/Actual	Cat	Dog	Goldfish
Cat	94	6	0
Dog	18	80	2
Goldfish	3	2	95

(a) Confusion matrix for the setup: Data aug

Predicted/Actual	Cat	Dog	Goldfish
Cat	94	6	0
Dog	9	85	6
Goldfish	1	1	98

(b) Confusion matrix for the setup: FuCiTNet, Data aug, FT, $\lambda = 0.005$

Table 7: Cat-vs-dog-vs-goldfish confusion matrices for fold #2 in 3FCV for the best reference model(a) and FuCiTNet(b)

Quantitative results. The performance, in terms of accuracy, of the best classification model with and without applying FuCiTNet, on the three datasets is presented in Tables 2, 4 and 6. Several configurations were analyzed, 'None' indicates that the network was initialized with ImageNet weights and retrained on the dataset. 'Data aug' indicates that the set of the optimal data augmentation techniques was applied. 'FT', for Fine Tuning, indicates that only the last fully connected layer was retrained on the dataset and the remaining layers were frozen. If 'FT' is not indicated, it means that the whole layers of the network were re-trained on the dataset. To compare FuCiTNet with the state-of-the-art, we also include the accuracy of the method proposed in

[7] on cat-vs-dog dataset only (see Table 2) and the method proposed in [9] on the three datasets (see Tables 2, 4 and 6).

When applying FuCiTNet to the test images, the classification model reaches 2.24%, 1.53% and 0.54% higher accuracy than the best classification model on cat-dog, church-vs-palace and cat-vs-dog-vs-goldfish respectively. The number of TP and TN also improves in all the three datasets as shown in the confusions matrices in Tables 3, 5 and 7.

With respect to the state-of-the-art, FuCiTNet provides 14%, 2% and 16.4% better accuracy than the approach proposed in [9] on cat-dog, church-vs-palace and cat-vs-dog-vs-goldfish respectively and 18.44% better accuracy than the approach proposed by [7] on cat-vs-dog.

These results can explain that FuCiTNet makes the object-class in the input images more distinguishable to the model. Our transformation approach can be considered as an incremental optimization to the rest of optimizations space, since on all the datasets, the best results were obtained by combining FuCiTNet with fine-tuning and data-augmentation.

We have also analyzed the mean confidence of each class in the three problems such that:

$$\text{mean confidence}(k\text{-class}) = \frac{\sum_{i=1}^T P_k(i)}{T}$$

where T is the number of correctly classified test images, from the test split, with k ground truth class. The mean confidence of the best reference model and FuCiTNet are shown in the third column of Tables 2, 4 and 6. As we can observe from these results FuCiTNet clearly improve the mean confidence of the model in all the classes in cat-vs-dog and church-vs-palace datasets. Although FuCiTNet have lowered the mean confidence in the cat-vs-dog-vs-goldfish dataset, the accuracy, number of true positives and true negatives have improved.

FuCiTNet provides several advantages over the method proposed in [7] and [9]. It does not require neither paired datasets of high and low resolution input images nor super-resolution network. Unlike [7], in which the transformations were learnt regardless the class of the sample, our results demonstrate that exploiting class-inherent features improves significantly the borders between visually similar classes.

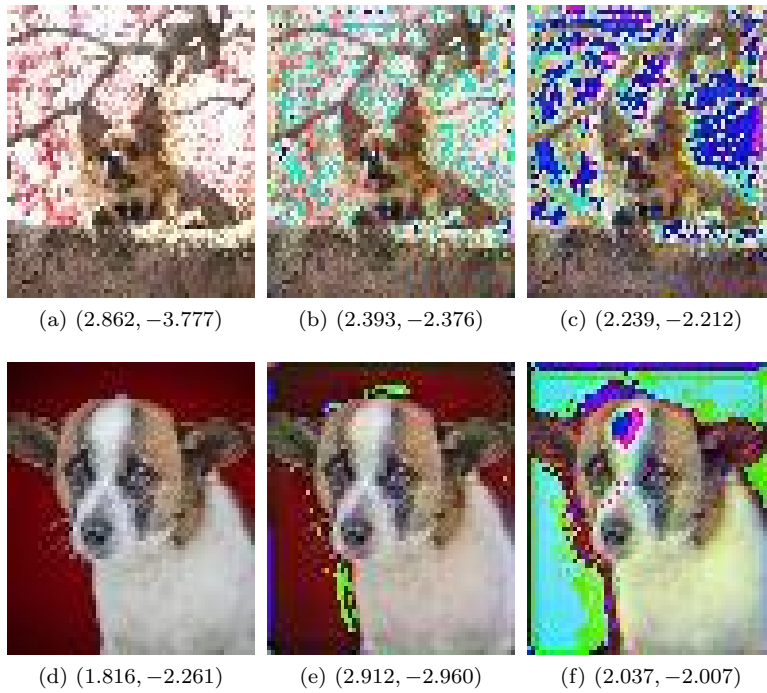


Figure 3: Left) Original sample of dog class; middle) Dog transformation to original sample; right) Cat transformation to original sample. The confidence score of the model for (dog class, cat class) are indicated below each image. The score for the left image is obtained using the best reference model whereas the ones for the middle and right images are obtained with FuCiTNet.

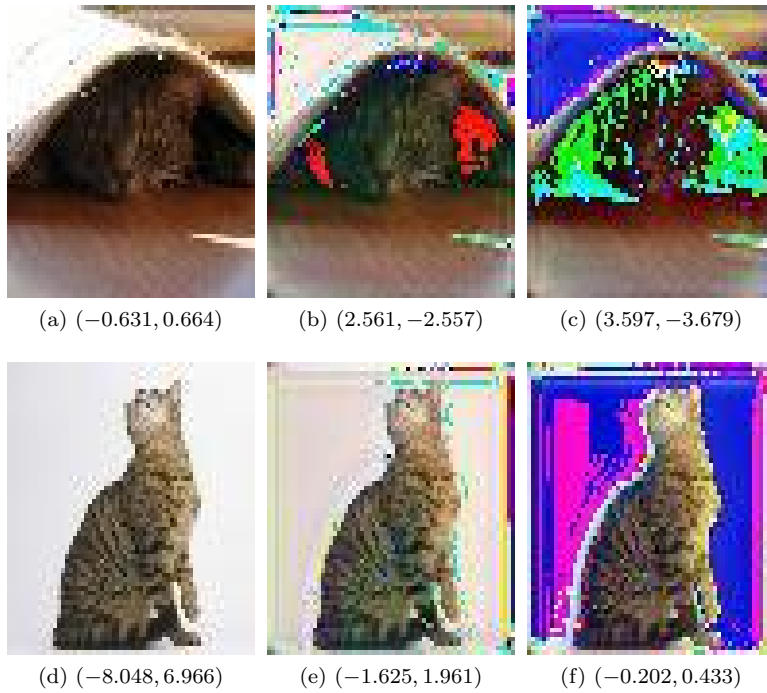


Figure 4: left) Original sample of cat class; middle) Dog transformation to original sample; right) Cat transformation to original sample. The confidence score of the model for (dog class, cat class) are indicated below each image. The score for the left image is obtained using the best reference model whereas the ones for the middle and right images are obtained with FuCiTNet.

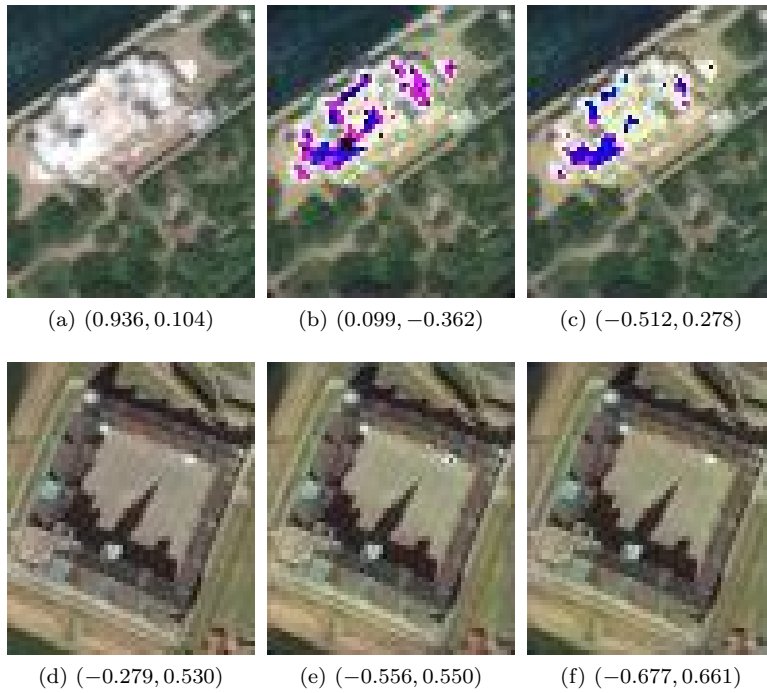


Figure 5: left) Original sample of Palace class; middle) Church transformation to original sample; right) Palace transformation to original sample. The confidence score of the model for (church class, palace class) are indicated below each image. The score for the left image is obtained using the best reference model whereas the ones for the middle and right images are obtained with FuCiTNet.

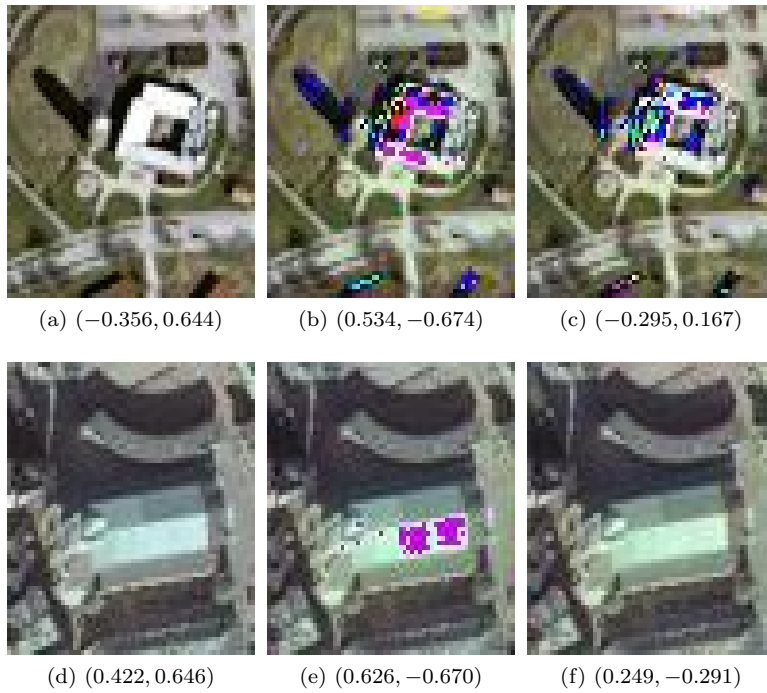
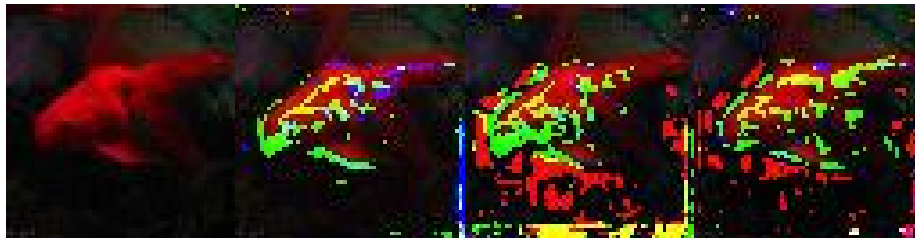


Figure 6: left) Original sample of Church class; middle) Church transformation to original sample; right) Palace transformation to original sample. The confidence score of the model for (church class, palace class) are indicated below each image. The score for the left image is obtained using the best reference model whereas the ones for the middle and right images are obtained with FuCiTNet.



(a) (-0.91, -2.78, 3.68) (b) (-6.06, -1.10, 5, 57) (c) (-6.71, -1.17, 5.99) (d) (-5.41, -1.11, 5.03)



(e) (-1.41, -1.56, 2.82) (f) (-3.05, -0.83, 3.77) (g) (-2.99, -0.80, 3.68) (h) (-3.12, -0.88, 3.91)

Figure 7: left) Original sample of goldfish class; middle left) Cat transformation to original sample; middle right) Dog transformation to original sample; right) Goldfish transformation to original sample. The confidence score of the model for (dog class, cat class, goldfish class) are indicated below each image. The score for the left image is obtained using the best reference model whereas the others are obtained with FuCiTNet.

Qualitative results. Figures 3, 4, 5, 6 and 7 show visually the transformations applied by FuCiTNet to different original input images from different classes. As it can be observed from these images, the transformations affect:

- the contour or border of the object class in the transformed images, as it can be seen in Figures 3(c) and (f) and Figures 4(c) and (f).
- the pixels that constitute the body of the object-class, as it is the case in Figure 3(f) and Figures 5(b) and (c), Figures 6(b), (c) and (e), Figures 7(c), (d), (f), (g) and (h).
- or background as it is the case in Figures 3(c) and (f), Figures 4(b), (c), (e) and (f).

In some images in the church-vs-palace testset, the model does not add any transformation. This indicates that the classifier does not need any additional transformation to differentiate the object-class in those specific images. In general, in this dataset, most of the palaces have either a dome or a rectangular roof with an inner courtyard but generally the churches do not include any of these features. In the cases where a church has an inner courtyard, the model needs to add a transformation so that the classifier can distinguish it better from a

church as it can be seen in Figure 6. Likewise, when the palace has a very distinguishable shape from a church, the model does not add any transformation as it can be seen in Figure 5(d), (e) and (f).

This effect does not occur in the cat-vs-dog and cat-vs-dog-vs-fish test sets, the background in these images is so diverse, it includes, people, furniture, occlusion etc., that the model always need some more transformation to differentiate better between the involved classes.

6. Conclusions and future work

Our aim in this work was reducing over-fitting in very small datasets in which the involved object-classes share many visual features. We presented FuCiTNet model in which a novel array of generators learn independently class-inherent transformations of each class of the problem. We introduced a classification loss in the generators to drive the learning of specific k-class transformations. The learnt transformations are then applied to the input test images to help the classifier distinguishing better between different classes.

Our experiments demonstrated that FuCiTNet increases the classification generalization capability on three small datasets with very similar classes. With the benchmark datasets we demonstrate that FuCiTNet behaves robustly in diverse-nature data and handles properly different view dimensions (zenital and frontal). We conclude that our method yields strong gains as an incremental optimization technique additional to the standards when searching for a better model performance.

As future work, we are planning to explore and adapt FuCiTNet to small medical datasets.

7. Acknowledgments

This work was partially supported by the Spanish Ministry of Science and Technology under the project TIN2017-89517-P and the project TEC2016-75976-R, financed by the Spanish Ministerio de Economía, Industria y Competitividad and the European Regional Development Fund (ERDF). S. Tabik was supported by the Ramon y Cajal Programme (RYC-2015-18136). E.G was supported by the European Research Council (ERC Grant agreement 647038 [BIODESERT]), with additional support from Generalitat Valenciana (CIDEGENT/2018/041).

References

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, F.-F. Li, Imagenet: A large-scale hierarchical image database, IEEE Conference on Computer Vision and Pattern Recognition (2009) 248–255.
- [2] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft COCO: common objects in context (2014). [arXiv:1405.0312](https://arxiv.org/abs/1405.0312).

- [3] A. Krizhevsky, Learning multiple layers of features from tiny images (2012).
- [4] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in Neural Information Processing Systems 25*, 2012, pp. 1097–1105.
- [5] P. Y. Simard, D. Steinkraus, J. C. Platt, Best practices for convolutional neural networks applied to visual document analysis, in: *Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 958–963.
- [6] H. S. Baird, Document image defect models, in: *Structured Document Image Analysis*, 1992, pp. 546–556.
- [7] L. Perez, J. Wang, The effectiveness of data augmentation in image classification using deep learning (2017). [arXiv:1712.04621](#).
- [8] X. Zhu, Y. Liu, Z. Qin, J. Li, Data augmentation in emotion classification using generative adversarial networks (2017). [arXiv:1711.00648](#).
- [9] A. J. Ratner, H. R. Ehrenberg, Z. Hussain, J. Dunnmon, C. Ré, Learning to compose domain-specific transformations for data augmentation (2017). [arXiv:1709.01643](#).
- [10] Y. Lecun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, *Proceedings of the IEEE* 86 (11) (1998) 2278–2324.
- [11] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: A simple way to prevent neural networks from overfitting, *Journal of Machine Learning Research* 15 (56) (2014) 1929–1958.
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: *Advances in Neural Information Processing Systems 27*, 2014, pp. 2672–2680.
- [13] B. Zoph, E. D. Cubuk, G. Ghiasi, T. Lin, J. Shlens, Q. V. Le, Learning data augmentation strategies for object detection (2019). [arXiv:1906.11172](#).
- [14] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition (2015). [arXiv:1512.03385](#).
- [15] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift (2015). [arXiv:1502.03167](#).
- [16] K. He, X. Zhang, S. Ren, J. Sun, Delving deep into rectifiers: Surpassing human-level performance on imagenet classification (2015). [arXiv:1502.01852](#).
- [17] V. Nair, G. Hinton, Rectified linear units improve restricted boltzmann machines vinod nair, in: *Proceedings of ICML*, Vol. 27, 2010, pp. 807–814.

- [18] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, W. Shi, Photo-realistic single image super-resolution using a generative adversarial network (2016). [arXiv:1609.04802](#).
- [19] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition (2014). [arXiv:1906.11172](#).
- [20] Y. Le, X. Yang, Tiny imagenet visual recognition challenge (2015).
URL <https://tiny-imagenet.herokuapp.com/>
- [21] G. Cheng, J. Han, X. Lu, Remote sensing image scene classification: Benchmark and state of the art, *Proceedings of the IEEE* 105 (10) (2017) 1865–1883.
- [22] D. Kingma, J. Ba, Adam: A method for stochastic optimization (2015). [arXiv:1412.6980](#).
- [23] A. Paszke, et al., Pytorch: An imperative style, high-performance deep learning library, in: *Advances in Neural Information Processing Systems* 32, 2019, pp. 8024–8035.