

# Disentanglement of Color and Shape Representations for Continual Learning

David Berga<sup>1</sup> Marc Masana<sup>1</sup> Joost Van de Weijer<sup>1</sup>

## Abstract

We hypothesize that disentangled feature representations suffer less from catastrophic forgetting. As a case study we perform explicit disentanglement of color and shape, by adjusting the network architecture. We tested classification accuracy and forgetting in a task-incremental setting with Oxford-102 Flowers dataset. We combine our method with Elastic Weight Consolidation, Learning without Forgetting, Synaptic Intelligence and Memory Aware Synapses, and show that feature disentanglement positively impacts continual learning performance.

## 1. Introduction

Convolutional Neural Networks have shown to increasingly achieve better performances in several recognition tasks over the past years (Krizhevsky et al., 2012; LeCun et al., 2015; Guo et al., 2016). In common image classification tasks, the network learns the whole dataset in a single training session. Thus, the network is only capable of doing inference on seen classes. If more classes would be learned without using a continual learning approach (i.e. finetuning on each task), the network would suffer of what is known as *catastrophic forgetting* (McCloskey and Cohen, 1989; French, 1999; Kirkpatrick et al., 2017). Catastrophic forgetting appears when neurons optimize their weights for a new task without taking into account previous knowledge; meaning previous classes performance is buried in favor of the new ones. Lifelong Learning and Continual Learning (CL) propose a more realistic scenario where the learner continually adapts to a sequence of tasks while avoiding said catastrophic forgetting. One of the reasons of catastrophic forgetting could be the lack of adaptability from network parameters in the stability-plasticity trade-off from task to task, this interference can cause the network to entangle all trained representations.

<sup>1</sup>Computer Vision Center, Bellaterra, Barcelona, Spain. Correspondence to: David Berga <dberga@cvc.uab.es>.

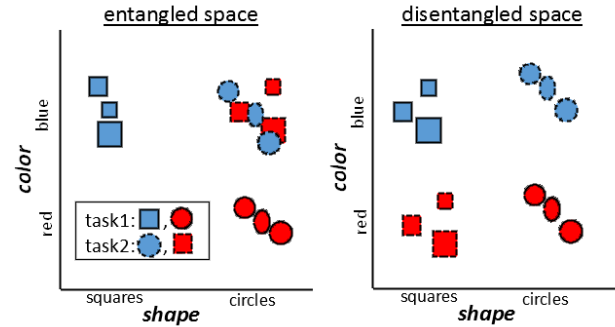


Figure 1. Illustration of feature disentanglement for a simple case of color "red/blue" and shape "square/circle". Left: Blue circles and red squares are entangled in the binding phase. Right: In a disentangled space, all shapes are separated along with colors.

Disentanglement of features, such as illumination, viewpoint object orientation or surface reflectance has been a long desired objective in computer vision (Tappen et al., 2003). It is also believed to play an important role in the success of deep learning (Bengio, 2009). In this paper, we hypothesize that disentanglement of features also plays an important role in continual learning settings. The main idea is that disentangled features can better generalize to new tasks (Bengio et al., 2013). As an example (Fig. 1), consider a system which should learn two tasks. The first task requires to distinguish between red circles and blue squares, while the second task requires to distinguish between blue circles and red squares. A network which would have learned a shape-color disentangled representation on the first task can easily adapt to the second task. However, a network which has learned an entangled representation (neurons firing for red-circles) might have more problems to generalize to other feature. While learning the second task, neurons specific to detect red-circles would instead detect blue-circles, consequently leading to catastrophic forgetting on the previous task.

Motivated by the biological evidence for separate processing of color and shape, we propose a two-branch network for image classification which fuses both at the end. This explicit color and shape disentanglement allows us to assess feature representation importance for continual learning and a way to fuse networks before binding.

## 2. Related Work

Latest reviews on CL (Parisi et al., 2019; Lange et al., 2019; Maltoni and Lomonaco, 2019) focus on evaluating approaches by equally distributing classes from a dataset into multiple tasks. Those tasks are then learned in an incremental fashion by fitting the network parameters to each new group of classes. To avoid catastrophic forgetting, approaches regularize the model, store information or replay data (Rebuffi et al., 2017; Shin et al., 2017; Lopez-Paz and Ranzato, 2017). Some of the first methods applied to neural networks, are focused on regularizing the weights or feature representations in order to keep those as close as possible to the older weights or representations while learning the new task. Learning without Forgetting (LwF) (Li and Hoiem, 2016) adds a regularization term to the cross-entropy loss which tries to push the outputs of previous classes to be similar to the outputs of new tasks before learning the task at hand. Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017) calculates the Fisher Information for all weights, which is then used as an importance measure on the regularization loss. We consider these approaches in addition to Synaptic Intelligence (SI) (Zenke et al., 2017) and Memory Aware Synapses (MAS) (Aljundi et al., 2018).

AlexNet (Krizhevsky et al., 2012; Flachot and Gegenfurtner, 2018) used 2 GPUs processing images in different groups of convolutions in parallel, sharing information in certain layers. By visualizing the filters in the first convolutional layer, the network showed sparse features that were distinct on the 2 network branches. These were similar to "gabor-like" filters for the gray branch and sinusoidal/concentric filters for the color branch, similar to receptive fields in V1 (see (Krizhevsky et al., 2012)-Fig. 3 and (Flachot and Gegenfurtner, 2018)-Fig. 2). Rafegas et al. (Rafegas and Vanrell, 2017; 2018; Rafegas et al., 2019) indexed selectivity of individual neurons to specific features in a VGG-M. By grouping these by color statistics (e.g. Hue), some appeared to be highly color selective and others low- or non color selective. Previous work (Khan et al., 2012) proposed an algorithm that processes shape and color separately (using SIFT features for shape features and color naming and pixel-wise hue descriptors for color features). Fusing these features showed improved accuracy in classification. From the aforementioned studies, we think a disentanglement procedure focused on separating color and shape feature computations in a network could be useful to acquire higher accuracy as well as preventing catastrophic interference, due to the specificity of neurons to each of these features.

## 3. Proposed Method

In this paper, we propose disentangled architecture able to prevent catastrophic inference restricting specific feature representations. We process color and Shape features

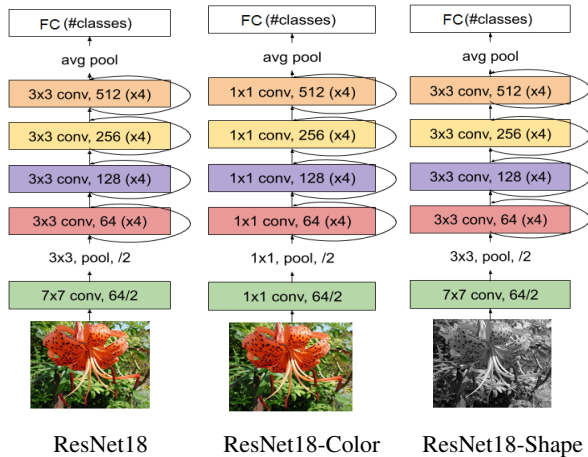


Figure 2. Structure of the standard ResNet18 network, color and shape networks. Every layer (represented in a different color) sums (+) its activation at every residual block.

separately. We do so using a ResNet18-DS architecture (Fig. 3) processing each feature representation separately in each branch (Fig. 2). The *objectives* are: a) Compute representations of Color and Shape separately. b) Learn disentangled representations for CL while retaining capacity. c) Provide an architectural mechanism able to be applied in combination with other CL algorithms (approach- and model-agnostic).

### 3.1. Independent Representations of Color and Shape

We first discuss the networks we use for the computation of shape-only, and color-only features. Separation of these two features might better represent task-dependent combinations of color and shape (see Fig. 1). We consider two network parts: the *Feature Extractor (FE)* and the *Head / Classifier*. Here we use a ResNet18 as our base network, (Fig. 2-Left), but the idea could be applied to other models. The network feature extractor is composed of a convolutional+maxpooling layer, followed by four convolutional blocks with batch normalization, and an average pooling layer, and the classifier contains a fully-connected (FC).

To uniquely process color (*ResNet18-Color*), we changed all convolutional operations (in ResNet18 are conv7x7 in first layer and conv3x3 in all blocks) by a convolution of 1x1 kernels (conv1x1). This *ResNet18-Color* network (Fig. 2-Mid) will learn each pixel independently from spatial computations, therefore it is only able to process color information. The absence of filters with spatial extend prevent it from being able to learn shape information.

For the case of shape (*ResNet18-Shape*), we transformed the RGB image to grayscale and used the original ResNet18 architecture. The *ResNet18-Shape* network (Fig. 2-Right) is processing local information only using intensity information (thus, unable to process RGB chromaticities).

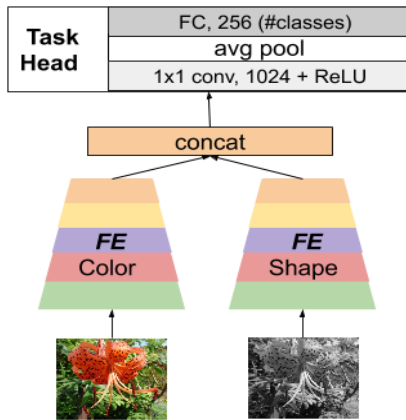


Figure 3. Disentangled network (ResNet18-DS) based on *Late Fusion* of ResNet18-Color and ResNet18-Shape. The 3 layers on the Head (gray) are independent for each task.

### 3.2. Fusing Color and Shape Representations in a Multi-Branch network (ResNet18-DS)

Having established architecture which process color and shape separately in the previous section, we here propose an architecture to combine the two branches. The main requirements of our architecture are:

- *Feature disentanglement*: The layers which are shared between the different tasks should only contain disentangled color and shape features. The entanglement should only happen in the task-specific head.
- *Spatial binding*: the combination of color and shape should be entangled before any pooling operation.

A system which would lack spatial binding is one where you would perform average pooling on the feature extractor of the color and shape branch and then combine the information. Such a system would know there are certain colors and shapes present, but could not say with certainty which color is connected to which shape. In Fig. 3 the architecture which fulfills our design requirements is presented. The two-branch network forwards two versions of the input (color and gray) to the FE part of each branch (color and shape). Next the output of both branches is concatenated (the output of each branch is 7x7x512 and after concatenation the dimension is 7x7x1024). Until here the color and shape information are disentangled and all layers are shared among the tasks. Then processing moves on to task-specific heads, in which the entanglement is performed. The task-specific head consists of four layers. First a 1x1 convolutional layer which maps the disentangled features, to entangled task specific features, followed by a ReLU. Then we perform an average pooling operation to the output of the convolution and flatten to a feature vector. Finally, a Linear layer (FC) maps the entangled features to the number of classes of each particular task. A softmax operation is added after the FC before computing the gradient loss.

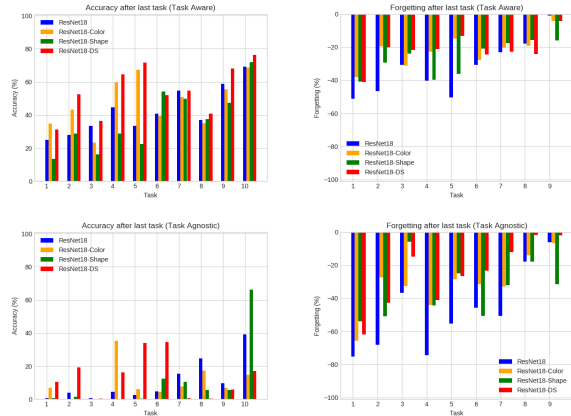


Figure 4. Average classification results (Left: Accuracy, Right: Forgetting) per task in finetuning on Oxford-102 dataset.

## 4. Experiments

We performed a set of experiments with the standard ResNet18 as well as with the ResNet18-Color, ResNet18-Shape and our ResNet18-DS (Color+Shape) with the Oxford-102 dataset in an incremental setting of 10 tasks. To define this setting, we split the dataset (102 categories) into 10 tasks and processed the same network through each task split through the FE. We considered backpropagating the Head part for each task separately (task-aware) or concatenated in one unique Head (task-agnostic). The capacity of the FE of ResNet-18 is 11.2M, while the capacity of Shape and Color branches FE are 11.2M and 1.5M respectively. The total capacity ResNet18-DS FE is 12.6M (about 112.5% with respect the ResNet-18). For the each head (Linear Head and Task Head) there is a distinct capacity per task, given that we compared performance from original ResNet18 with same Task Head (ResNet18-H).

We processed distinct hyperparameters for each model to convergence (with learning rates 5e-2, 5e-3 and 5e-4 and batch size 32) for 200 epoch. We also applied a weight decay of 0.0002, momentum of 0.9 and we considered a patience of 15 (lowering the learning rates by a factor of 3 when loss does not improve after 15 epoch, until  $lr < 1e-6$ ).

In Table 1 we show results from the ablation results of each separate branch (ResNet18-Color and ResNet18-Shape). We have also processed the standard ResNet18 and our ResNet18-DS (Color+Shape) network with CL regularization algorithms (LwF, EWC, SI and MAS). Our model outperforms the standard state-of-the-art both in finetuning and with these two algorithms. Considering the classification results, our model presents lower forgetting with respect the standard ResNet, this means that our models is capable of disentangling activity of each task split at a feature level (the FE). In terms of accuracy per task our model performance is able to retain similar accuracy over tasks, whereas the standard ResNet18 accuracy lowers in both task-

## Disentanglement of Color and Shape Representations for Continual Learning

Task-Aware	ResNet18	ResNet18-H*	ResNet18-Color	ResNet18-Shape	ResNet18-DS
Finetune	42.4	40.1	47.8	37.0	54.8
LwF	58.3	55.8	53.8	50.9	<b>62.8</b>
EWC	46.4	45.4	52.9	41.3	55.5
MAS	52.0	43.6	50.1	43.0	56.9
SI	43.4	45.3	50.1	42.0	53.8
Task-Agnostic	ResNet18	ResNet18-H*	ResNet18-Color	ResNet18-Shape	ResNet18-DS
Finetune	10.6	10.2	10.0	10.2	13.8
LwF	13.7	9.8	7.4	11.6	12.6
EWC	10.3	11.3	12.4	9.1	<b>17.0</b>
MAS	11.2	7.1	9.0	8.3	11.8
SI	10.7	10.5	8.8	10.0	14.6

Table 1. Average accuracy classification results for Oxford-102 dataset in a incremental setting of 10 tasks. Mean Accuracy calculated from all tasks after the 10th task. \* Same task-specific head (and head capacity) as ResNet18-DS (Fig.3). **Bold** is TOP-1 accuracy.

Task-Aware	ResNet18	ResNet18-H*	ResNet18-Color	ResNet18-Shape	ResNet18-DS
Finetune	-29.2	-29.8	-19.8	-24.0	-19.3
LwF	-12.5	-13.4	-13.8	-14.9	-11.0
EWC	-22.6	-24.9	-14.3	-20.2	-17.1
MAS	-11.9	-9.1	-12.1	<b>-4.6</b>	-9.4
SI	-26.1	-25.3	-14.7	-19.4	-18.0
Task-Agnostic	ResNet18	ResNet18-H*	ResNet18-Color	ResNet18-Shape	ResNet18-DS
Finetune	-43.0	-47.6	-28.4	-31.3	-22.7
LwF	-31.6	-35.6	-24.1	-38.7	-23.5
EWC	-31.6	-38.6	-18.0	-23.6	-36.2
MAS	-9.5	-6.9	-8.3	<b>-2.7</b>	-9.3
SI	-34.9	-44.6	-20.2	-32.6	-30.7

Table 2. Average forgetting for Oxford-102 dataset in a incremental setting of 10 tasks. Mean Forgetting calculated as drop of accuracy from all tasks after the 10th task. \* Same task-specific head (and head capacity) as ResNet18-DS (Fig.3). **Bold** is lowest forgetting.

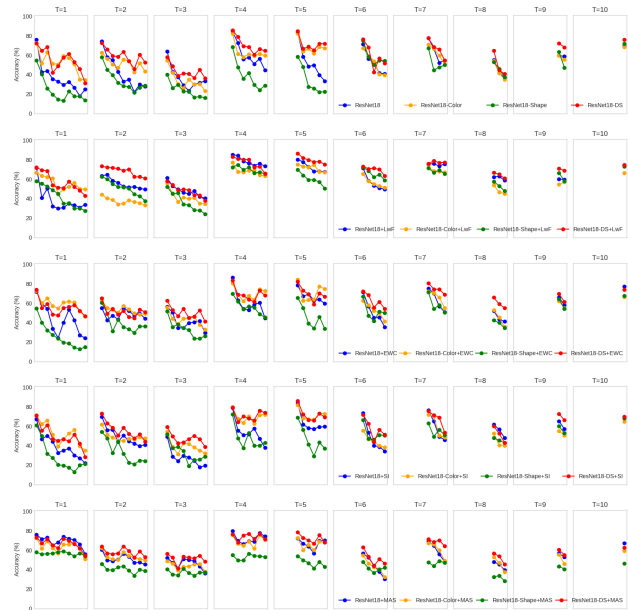


Figure 5. Task-sequence results for 10 tasks in Task-Aware setting for Finetune, LwF, EWC, SI and MAS.

aware and task-agnostic settings. We showed in Figs. 4-5 that our model acquires overall higher accuracy and lower forgetting across all tasks. We believe this is due to the interference between features of the FE (specially at latter

layers, which bind higher-level information). The standard ResNet shares all representations (e.g. color and shape) in a unique branch, which struggles on adapting weights for each new task representations. This interference prevented by our disentanglement procedure. We would like to point out that overall classification results for Task-Agnostic would be higher with exemplars.

## 5. Conclusion

In this study we propose a novel architectural design that allows disentanglement of color and shape representations in a convolutional neural network. This method prevents catastrophic interference between these feature types, showing lower forgetting in comparison with standard networks. We show this strategy can be useful to be combined with other CL approaches (outperforming results based on standard architectures) and the potential to be used with distinct architectures and datasets (e.g. shapes, attributes). As such, we hope this paper inspires more research into the importance of disentangled representation for continual learning.

## Acknowledgements

We acknowledge the support from Huawei Kirin Solution. Marc Masana acknowledges 2019-FI.B2-00189 grant from Generalitat de Catalunya.

## References

- Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Computer Vision – ECCV 2018*, pages 144–161. 2018.
- Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, August 2013.
- Yoshua Bengio. Learning deep architectures for ai. *Foundations and trends in Machine Learning*, 2(1):1–127, 2009.
- Alban Flachot and Karl R. Gegenfurtner. Processing of chromatic information in a deep convolutional neural network. *Journal of the Optical Society of America A*, 35(4):B334, March 2018.
- Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, April 1999.
- Yanming Guo, Yu Liu, Ard Oerlemans, Songyang Lao, Song Wu, and Michael S. Lew. Deep learning for visual understanding: A review. *Neurocomputing*, 187:27–48, April 2016.
- Fahad Shahbaz Khan, Joost Van de Weijer, and Maria Vanrell. Modulating shape features by color attention for object recognition. *International Journal of Computer Vision*, 98(1):49–64, 2012.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, March 2017.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25*, pages 1097–1105. 2012.
- Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. Continual learning: A comparative study on how to defy forgetting in classification tasks. *arXiv preprint*, arXiv:1909.08383, 2019.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. In *Computer Vision – ECCV 2016*, pages 614–629. 2016.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems*, pages 6467–6476, 2017.
- Davide Maltoni and Vincenzo Lomonaco. Continuous learning in single-incremental-task scenarios. *Neural Networks*, 116:56–73, August 2019.
- Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of Learning and Motivation*, pages 109–165. Elsevier, 1989.
- German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, May 2019.
- Ivet Rafegas and Maria Vanrell. Color representation in CNNs: Parallelisms with biological vision. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*. IEEE, October 2017.
- Ivet Rafegas and Maria Vanrell. Color encoding in biologically-inspired convolutional neural networks. *Vision Research*, 151:7–17, October 2018.
- Ivet Rafegas, Maria Vanrell, Luís A. Alexandre, and Guillem Arias. Understanding trained CNNs by indexing neuron selectivity. *Pattern Recognition Letters*, October 2019.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. iCaRL: Incremental classifier and representation learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, July 2017.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *Advances in Neural Information Processing Systems*, pages 2990–2999, 2017.
- Marshall F Tappen, William T Freeman, and Edward H Adelson. Recovering intrinsic images from a single image. In *Advances in neural information processing systems*, pages 1367–1374, 2003.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML17*, page 39873995, 2017.