# Time for a Background Check! Uncovering the impact of Background Features on Deep Neural Networks

**Vikash Sehwag** [1]  **Rajvardhan Oak** [2]  **Mung Chiang** [3]  **Prateek Mittal** [1]

## Abstract

With increasing expressive power, deep neural networks have significantly improved the state-of-the-art on image classification datasets, such as ImageNet. In this paper, we investigate to what extent the increasing performance of deep neural networks is impacted by background features? In particular, we focus on *background invariance*, i.e., accuracy unaffected by switching background features and *background influence*, i.e., predictive power of background features itself when foreground is masked. We perform experiments with 32 different neural networks ranging from small-size networks (Howard et al., 2019) to large-scale networks trained with up to one Billion images (Yalniz et al., 2019). Our investigations reveal that increasing expressive power of DNNs leads to higher influence of background features, while simultaneously, increases their ability to make the correct prediction when background features are removed or replaced with a randomly selected texture-based background.

## 1. Introduction

One key driver behind this success of modern deep neural networks (DNNs) is their expressive power, which enables them to learn a rich set of representations required to solve a target task (Krizhevsky et al., 2012; Bengio et al., 2013; He et al., 2016). However, this expressive power raises a number of fundamental questions: how good are the learned representations? For example, does the learned representations correspond to semantic features in the image? Is the representation overfitted to the training data, or is it robust to the addition of semantically unrelated features in the image?

In this paper, we argue that one fundamental principle for

representation learning should be that the network should achieve background invariance (Pinto et al., 2008) i.e., be able to make a correct prediction when the background features are switched in an image classification task. Thus, we aim to understand to what extent current deep neural networks utilize the background information itself to solve the task of image classification.

We develop a framework where we examine the influence of background image features on output prediction using two approaches: 1) We mask the foreground content of all test images, thus examine the correlation of output label with background features, 2) we switch the background of each image with an artificial background, such as white or texture-based images. In this setup, we expect the accuracy to remain unchanged assuming that the network has learned representations that are invariant to background changes. We rigorously test this framework on the ImageNet (Deng et al., 2009) dataset where we evaluate the performance of 32 different DNNs proposed over the last couple of years (Canziani et al., 2016).

Our findings are intriguing: we observe that even on a diverse dataset like ImageNet, state-of-the-art DNNs ends up learning correlation with background signal. In fact, even when we mask the foreground, we find that DNNs can make a correct prediction for a significant number of images. We highlight a few such examples in Figure 1, where just with the background content, a ResNet-18 network is able to make the correct prediction. We discuss this observation in more detail in Section 3. Similarly, when we switch image backgrounds while keeping the foreground unchanged, we observe a significant reduction in the test accuracy. It further supports our hypothesis that current DNNs do really exploit background correlation to make the correct prediction. Although, with increasing expressive power, the drop in accuracy after a random background switch keeps decreasing.

Since the current training loss functions, such as cross-entropy loss (Goodfellow et al., 2016), do not explicitly encourage background invariance, we find that increasing the expressive power of current networks further enables them to exploit any existing correlation between the background features and output prediction (supporting detail

---

[1]Princeton Universty, USA [2]University of California, Berkeley, USA [3]Purdue University, USA. Correspondence to: Vikash Sehwag <vvikash@princeton.edu>.
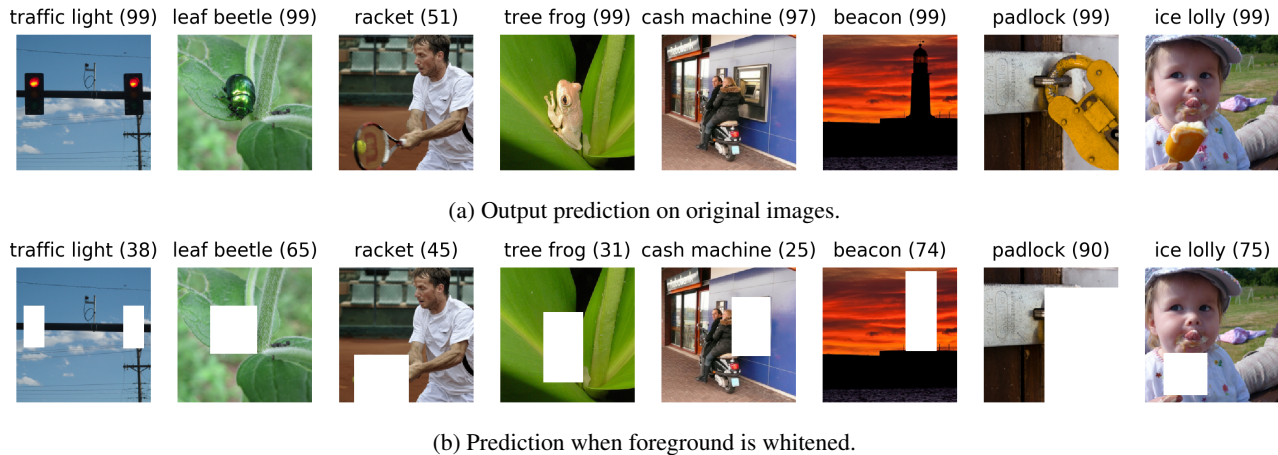
traffic light (99)   leaf beetle (99)   racket (51)   tree frog (99)  cash machine (97)  beacon (99)   padlock (99)   ice lolly (99)

(a) Output prediction on original images.

traffic light (38)   leaf beetle (65)   racket (45)   tree frog (31)  cash machine (25)  beacon (74)   padlock (90)   ice lolly (75)

(b) Prediction when foreground is whitened.

*Figure 1.* Predicted label of original and background only images (output confidence in parenthesis) for a ResNet-18 network trained on ImageNet dataset. It shows that even when the foreground is absent, the network is able to make a correct prediction for a large number of images. Such correlation with background features can a sign of 1) a network utilizing context from background to make correct prediction or 2) a network exploiting the inherent bias in the data, thus using background correlation itself to discriminate between classes.

in Section 3). Our results further question the ability of the deep neural networks to learn fundamental semantics, necessary to solve the task at hand, in the *current training paradigms*. While further increasing the dataset size or diversity, thus training on a larger set of background variations, could be one approach to increase background invariance, we believe that a more successful approach will be to improve the training loss function to penalize correlation with the background.

**Contributions.** We present a rigorous evaluation of 32 different DNNs to demonstrate that output prediction of deep neural networks is heavily influenced by background features. Our first set of results shows that with increasing expressive power, the tendency of DNNs to exploit background features to solve the task at hand increases. Next, we demonstrate a large reduction in the performance of DNNS when we switch the background features across multiple datasets and network architectures. However, this gap in performance decreases with increasing expressive power of DNNs.

## 2. Methods and experimental setup

Most computer vision datasets, such as ImageNet, VOC12, Caltech101 have object categories corresponding to a noun. This structure enables us to divide each image into two parts: Foreground, which comprises the correct objects, and background, which is everything in the image except the foreground. Note that for datasets based on scenes (SUN (Xiao et al., 2010)), texture (DTD (Cimpoi et al., 2014)) such categorization of foreground and background is not feasible.

Our objective is to analyze the following question: To what

extent current deep neural networks exploit background features to achieve the targeted classification. We approach this question from two directions.

**Testing background influence by masking foreground.** We first aim to explicitly measure the correlation between background features and output prediction. We mask the foreground, i.e., the object corresponding to correct prediction and replace it with a white patch. We choose a white patch, instead of random noise to have a minimal influence of the patch on output. We refer to such images as background images.

**Testing background invariance by switching background.** Now we aim to analyze the effect of the absence of background features. We argue that if the learned representations are agnostic to the background features, replacing it with a white image or other novel patterns like texture, outdoor scenes won't affect it. Thus in this setup, we expect networks to achieve accuracy close to unmodified images.

We also focus on relative accuracy, instead of an absolute number for accuracy, for both of the aforementioned tasks. We argue that it potentially allows us the analyze the performance on images unaffected by labeling noise in the foreground and background classification. We discuss this choice in detail in Appendix A.

**Experimental setup.** We use 31,801 images from the validation set of ImageNet dataset (data preparation process described in Appendix B). We evaluate 32 different networks[1] in our framework. These networks bring in variation in depth (ResNet-18 to resnet101), variation in architec-

---

[1] we use publicly available checkpoints from https://github.com/rwightman/pytorch-image-models.

*Table 1.* Analyzing the correlation between the expressive power of a neural network and its tendency to latch on background information to make correct prediction. ✓and ✗ implies whether the background features are itself classified as the true class or not. It shows that with increasing expressive power, we found a new set of images, where the background features are sufficient to achieve correct classification for networks with equal or higher expressive power, but not for networks with lower expressive power.

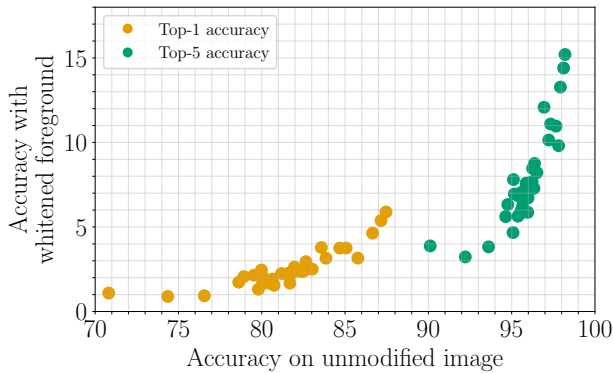| Architecture | Accuracy on unmodified images (%) | Accuracy with whitened foreground (%) | racket | scoreboard | grand piano | dumbbell | custard apple | missile | mouse |
|---|---|---|---|---|---|---|---|---|---|
| ResNet-18 (He et al., 2016) | 74.3 | 0.89 | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| MobileNet-v3 (Howard et al., 2019) | 76.5 | 0.93 | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| ResNet101 (He et al., 2016) | 80.7 | 1.55 | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Wide-ResNet-50-2 (Zagoruyko & Komodakis, 2016) | 81.7 | 2.28 | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| EfficientNet-b1 (Tan & Le, 2019) | 83.6 | 3.78 | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ |
| Swsl-ResNext50-32x4d (Yalniz et al., 2019) | 84.6 | 3.75 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |
| ResNext101-32x48d (Mahajan et al., 2018) | 87.4 | 5.88 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |



*Figure 2.* Top-1 and top-5 accuracy for background features when the foreground content is masked. Each dot corresponds to one neural network. It shows that with increasing expressive power, current DNNs increasingly exploit correlation of background features and correct prediction.

ture (ResNet, DPN, MobileNet, etc) (He et al., 2016; Chen et al., 2017; Howard et al., 2019)), and effect of training setup (models trained on ∼ 1 Billion images using the semi-weakly supervised training (Yalniz et al., 2019; Mahajan et al., 2018)). We primarily use white and texture-based images as two alternatives to replace the image background. We demonstrate a similar trend with other choices such as background based on Gaussian or uniform noise or any random color. Similarly, we also show a similar trend as ImageNet for Caltech101 and VOC12 dataset. We use the ResNet18 network as the baseline network.

## 3. Background Influence: DNNs learns correlation between correct labels and background features

Now we present the results with foreground masking, where we mask the foreground and test whether deep neural networks have learned correlation with background features.

Figure 2 presents these results where we measure both top-1 and top-5 accuracy of 32 different networks on only back-

ground features. We plot these results in order of their test accuracy on unmodified ImageNet images and observe that the influence of background features increases with the expressive power of DNNs. Such correlation with background features can a sign of 1) either a network utilizing context from background to make a correct prediction or 2) or a network exploiting the inherent bias in the data, thus using background correlation itself to discriminate between classes.

We analyze this phenomenon in more detail in Table 1. We consider seven networks with different expressive power and report a set of images where they make correction predictions solely based on background features. It shows that there is a common set of images where all networks are able to make a correct prediction with background features. However, as we increase expressive power, we find a novel subset of images for which background features are not predictive for any less expressive networks but able to make the right prediction for all more expressive networks.

## 4. Background invariance: Switching background leads to performance degradation

Here we discuss the results with the background switching experiment where we replace the original background features with white or texture-based images.

Figure 4 presents the top-1 accuracy for all 32 networks evaluated in this experiment (top-5 accuracy in Appendix C). If the networks had learned a set of good representation, the accuracy should not change with a change in the background. However, our results show a large drop in performance when we switch background features. However, with increase expressive power of DNNs, we observe that this gap decreases. We present a few examples demonstrating this phenomenon in Figure 3.

When tested with the baseline ResNet-18 network, we observe a similar drop in performance for Gaussian noise, Uni-

(a) Using white background

(b) Using texture-based background

*Figure 3.* Mis-classification for ResNet-18 network when we change the background features.

form noise, green, yellow, grey color-based backgrounds. Similarly, for the VOC12 dataset (Everingham et al., 2015), test accuracy decreases to 75% and 46% when using white and texture-based backgrounds, respectively. Original images achieved 95% test accuracy on the same dataset. For the Caltech101 dataset (Fei-Fei et al., 2004; 2006), test accuracy decrease to 85% and 74%, when using white and texture-based backgrounds, respectively. We report the mean accuracy over 3 runs. Original images achieved 94% test accuracy on the same dataset.
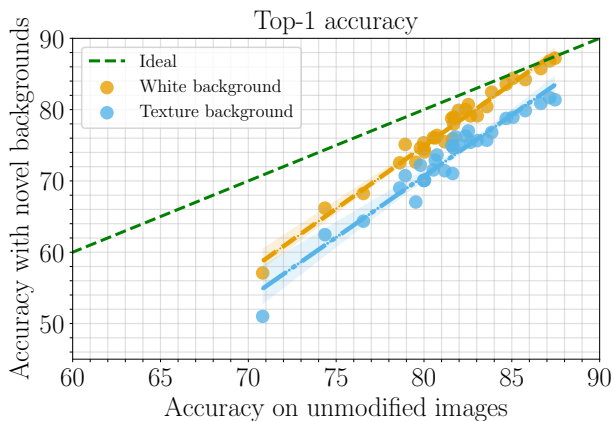


*Figure 4.* Top-1 accuracy on ImageNet test images when the background is switched to white or texture based background. Each dot correspond to one neural network.

## 5. Related work

Our work is closely related to an existing line of research (Zhu et al., 2016; Torralba, 2003; Sun & Jacobs, 2017; Katti et al., 2019; Geirhos et al., 2020), and a concurrent work (Xiao et al., 2020), aiming to understand the relationship of foreground and background features with network performance. However, our work particularly delves deeper into understanding this relationship with the increasing expressive power of DNNs. Another related line of research in biological vision studies the importance of context (Oliva & Torralba, 2007; Hock et al., 1974; Biederman et al., 1982; Hayes et al., 2007) and other studies the impact of changing backgrounds (DiCarlo et al., 2012; Yamins et al., 2014; Pinto et al., 2008; George & Hawkins, 2009).

## 6. Discussion

In this work, we demonstrate a significant impact of background features on the performance of deep neural networks. We provide supporting experimental results in two frameworks: Evaluating performance 1) only on background features 2) while removing or switching background features. With increasing expressive power of DNNs, we observe an increasing tendency to exploit background features to get the correct prediction, while simultaneously increasing the ability to make a correct prediction with foreground features only. Future works should aim to delve deeper into the correlation with background features while also measuring the performance with the worst-case background, instead of a random background.

# References

Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.

Biederman, I., Mezzanotte, R. J., and Rabinowitz, J. C. Scene perception: Detecting and judging objects undergoing relational violations. *Cognitive psychology*, 14(2): 143–177, 1982.

Canziani, A., Paszke, A., and Culurciello, E. An analysis of deep neural network models for practical applications. *arXiv preprint arXiv:1605.07678*, 2016.

Chen, Y., Li, J., Xiao, H., Jin, X., Yan, S., and Feng, J. Dual path networks. In *Advances in neural information processing systems*, pp. 4467–4475, 2017.

Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., , and Vedaldi, A. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

DiCarlo, J. J., Zoccolan, D., and Rust, N. C. How does the brain solve visual object recognition? *Neuron*, 73(3): 415–434, 2012.

Everingham, M., Eslami, S. A., Van Gool, L., Williams, C. K., Winn, J., and Zisserman, A. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.

Fei-Fei, L., Fergus, R., and Perona, P. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, pp. 178–178. IEEE, 2004.

Fei-Fei, L., Fergus, R., and Perona, P. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.

Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Shortcut learning in deep neural networks. *arXiv preprint arXiv:2004.07780*, 2020.

George, D. and Hawkins, J. Towards a mathematical theory of cortical micro-circuits. *PLoS computational biology*, 5 (10), 2009.

Goodfellow, I. J., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, Cambridge, MA, USA, 2016. http://www.deeplearningbook.org.

Hayes, S. M., Nadel, L., and Ryan, L. The effect of scene context on episodic object recognition: parahippocampal cortex mediates memory encoding and retrieval success. *Hippocampus*, 17(9):873–889, 2007.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

Hock, H. S., Gordon, G. P., and Whitehurst, R. Contextual relations: the influence of familiarity, physical plausibility, and belongingness. *Perception & Psychophysics*, 16 (1):4–8, 1974.

Howard, A., Sandler, M., Chu, G., Chen, L., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al. Searching for mobilenetv3. arxiv 2019. *arXiv preprint arXiv:1905.02244*, 2019.

Katti, H., Peelen, M. V., and Arun, S. Machine vision benefits from human contextual expectations. *Scientific reports*, 9(1):1–12, 2019.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.

Mahajan, D., Girshick, R., Ramanathan, V., He, K., Paluri, M., Li, Y., Bharambe, A., and van der Maaten, L. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 181–196, 2018.

Oliva, A. and Torralba, A. The role of context in object recognition. *Trends in cognitive sciences*, 11(12):520–527, 2007.

Pinto, N., Cox, D. D., and DiCarlo, J. J. Why is real-world visual object recognition hard? *PLoS computational biology*, 4(1), 2008.

Sun, J. and Jacobs, D. W. Seeing what is not there: Learning context to determine where objects are missing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5716–5724, 2017.

Tan, M. and Le, Q. V. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.

Torralba, A. Contextual priming for object detection. *International journal of computer vision*, 53(2):169–191, 2003.

Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., and Torralba, A. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 3485–3492. IEEE, 2010.

Xiao, K., Engstrom, L., Ilyas, A., and Madry, A. Noise or signal: The role of image backgrounds in object recognition. *arXiv preprint arXiv:2006.09994*, 2020.

Yalniz, I. Z., Jégou, H., Chen, K., Paluri, M., and Mahajan, D. Billion-scale semi-supervised learning for image classification. *arXiv preprint arXiv:1905.00546*, 2019.

Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., and DiCarlo, J. J. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*, 111(23):8619–8624, 2014.

Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

Zhu, Z., Xie, L., and Yuille, A. L. Object recognition with and without objects. *arXiv preprint arXiv:1611.06596*, 2016.

## A. Setting up a baseline

ImageNet, the dataset we use in this work, does include some label noise in bounding boxes for background and foreground classification. For some images, we find that not all foreground content is labeled as foreground, e.g., not all ants in an image are labeled as foreground, thus even after removing foreground, the background image still has traces of the correct object. In that case, it is expected to obtain correct classification even with the background image. Similarly, if a significant part of the foreground image is dropped due to poor labeling, the network might not be able to achieve the correct classification with a foreground image.

To address this challenge, we focus on relative performance in our framework w.r.t. to a baseline network. We assume that any error on bounding boxes for foreground objects will affect the performance of all networks. Thus relative performance w.r.t. to a baseline network, potentially allows us the analyze the performance on images unaffected by labeling noise.

## B. Data Preparation

We use the images from the validation set of ImageNet dataset (50,000 images) following a two-step filtering process. First, we remove images with foreground area less than 5% of the image (for classes like ping-pong ball). With such images, when we mask the foreground, most networks still make the correct prediction. Similarly, with only foreground, most networks fail to make a correct prediction.

Next, we remove images where the foreground area is more than 70% of the image. We aim to understand the impact of background features where it desirable to have a significant ratio of background features in the image. Otherwise, if the image foreground covers the whole image, it is unlikely that background switching will cause misclassification.

Following this filtering process, we arrive at 31,801 images which we use throughout our experiments in this paper.

## C. Additional results for background switching experiment

Figure 5 presents the top-5 accuracy for all 32 networks evaluated in this paper. Similar to the top-1 accuracy, our results show a large drop in performance when we switch background features. However, with the increasing expressive power of networks, this drop in performance decreases.
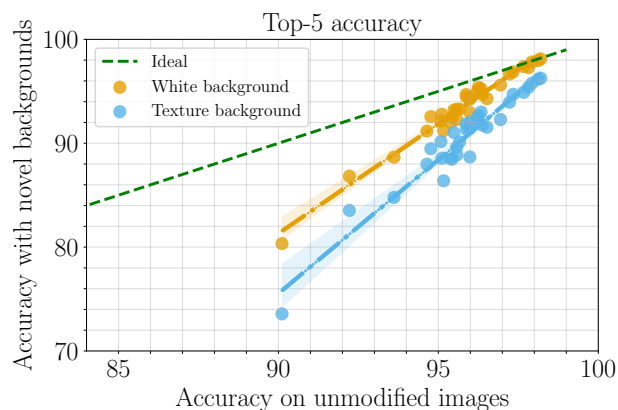


*Figure 5.* Top-5 accuracy on ImageNet test images when the background is switches to white or texture based background. Each dot correspond to one neural network.